000

Conformal Prediction for Hierarchical Data

Anonymous Authors¹

Abstract

We consider conformal prediction of multivariate data series, which consists of outputting prediction regions based on empirical quantiles of pointestimate errors. We actually consider hierarchical multivariate data series, for which some components are linear combinations of others. The intuition is that the hierarchical structure may be leveraged to improve the prediction regions in terms of their sizes for given coverage levels. We implement this intuition by including a projection step (also called reconciliation step) in the split conformal prediction [SCP] procedure and prove that the resulting prediction regions are indeed globally smaller than without the projection step. The associated strategies and their analyses rely on the literatures of both SCP and forecast reconciliation. We also illustrate the theoretical findings, both on artificial and on real data.

1. Introduction

This article combines two post-hoc procedures (two procedures that are applied after initial forecasts were computed): conformal prediction and forecast reconciliation for hierarchical data, both in a regression setting. Hierarchical data refers to multivariate observations abiding by some linear structure such that some components (referred to as in the higher levels of the hierarchy) are given by linear combinations of a subset of the components (referred to as the lower level of the hierarchy).

Multivariate conformation prediction. Conformal prediction is a general approach to output prediction sets, based on finite samples and on any underlying forecasting method, under mild assumptions—typically, exchangeability of data. It was first made formal by Vovk et al. (2005) and gained attention since the work by Lei et al. (2018). We are interested in the multivariate (also called multitarget) extensions of conformal predictions, discussed by Johnstone & Cox (2021), Messoudi et al. (2021; 2022), and Feldman et al. (2023); see Appendix E for more details. These works deal with prediction regions providing joint coverage guarantees, while we will rather be interested in component-wise coverage guarantees.

Forecast reconciliation. Forecast reconciliation is about taking into account the hierarchical structure of the multivariate data to improve forecasts. The intuition guiding this approach is that observations at the higher levels of the hierarchy are often easier to forecast, and that these forecasts can be leveraged to improve the forecasts at the lower levels of the hierarchy. Conversely, valuable local information from the lower levels can be leveraged to improve forecasts at the higher ones. We mainly focus on a series of works (Hyndman et al., 2011, Wickramasuriya et al., 2019, Panagiotelis et al., 2021) that approach reconciliation by projections onto the subspace of so-called coherent forecasts; see Appendix D for more details.

Contributions and challenges. This article combines (for the first time) both theories, to provide improved conformal predictions in the case of hierarchical data. The benchmark procedure consists of the split conformal prediction procedure (Lei et al., 2018) run on signed non-conformity scores (as in Linusson et al., 2014) and in a component-wise fashion. The improved version only differs from this benchmark through the application of a projection to the regression function learned on the train set-in the spirit of forecast reconciliation. However, the challenge, and our contributions, lie in showing that for a given coverage level, the resulting prediction regions indeed leverage the hierarchical structure of the data and are more efficient, i.e., smaller in some sense to be made precise. Actually (and only because we are interested in component-wise coverage, see Appendix E.2), the main difficulty was to state and control the corresponding criterion, referred to as $(\star\star)$ in the sequel. To do so, we show that one can resort to some known trace inequalities of the literature of forecast reconciliation, and also introduce some new such trace inequalities.

A more precise description of the challenges overcome may be found after the sketch of the proof of Theorem 3.7 in Section 3.2 (see also the specific discussions in Appendices E.2 and D.2).

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Unrelated references. The terminology "hierarchical" also appears in Lee et al. (2023), Dunn et al. (2023) and Duchi et al. (2024) in the contexts of predictive inference or con-058 formal prediction. However, this line of research is about 059 a completely different setting (data coming from different 060 sources) and bears no relationship to the hierarchical struc-061 ture we consider. Actually, Duchi et al. (2024) rather use the 062 terminology "multi-environment" in that context to avoid 063 potential confusions.

064

088

089

090

091

092

093

094

095

096

097

098

099

100

104

065 **Outline.** In Section 2, we formally state the setting consid-066 ered, the objectives targeted, and the methodology followed. 067 The objectives consist of a component-wise coverage objec-068 tive (\star) and of an efficiency (small-length) objective $(\star\star)$. 069 The methodology consists of taking the split conformal pro-070 cedure (run component-wise with signed non-conformity scores) as a benchmark, and discuss how to improve it by adding a reconciliation step through projections. Section 3 then states the theoretical results achieved: the coverage 074 guarantees in Theorem 3.2, and efficiency results (weak and 075 practical version in Theorem 3.7, strong and oracle version in Theorem 3.10). We only provide a sketch of the proof 077 of Theorem 3.7 (to give an idea of how we connected the 078 tools of conformal prediction and of forecast reconciliation), 079 and defer full proofs of all these results in Appendices A-Appendix B-Appendix C. Finally, Section 4 illustrates the 081 theoretical finds both on artificial and on real data, with full 082 details of the simulations to be found in Appendix F. The 083 real data concern the charging of electric vehicles. We recall that Appendices D and E discuss the literature of forecast 085 reconciliation and multivariate split conformal prediction, respectively, and review their classic results. 087

Notation. For an integer $n \ge 1$, let $[n] = \{1, \ldots, n\}$. For a real number $x \ge 0$, we let |x| and $\lceil x \rceil$ denote the lower and upper integer parts, respectively. For a vector $oldsymbol{u} \in \mathbb{R}^m$ and $n \leq m$, let $\boldsymbol{u}_{1:n} = (u_1, \dots, u_n)^{\mathsf{T}}$ be the vector of the first *n* components of \boldsymbol{u} . The null vector of \mathbb{R}^m is denoted by $\mathbf{0} = (0, \dots, 0)^{\mathsf{T}}$. We let diag (\boldsymbol{w}) denote the $m \times m$ diagonal matrix with diagonal elements given by $w \in \mathbb{R}^m$. We denote by Id_m the $m \times m$ identity matrix. The trace of square matrix M is denoted by Tr(M).

2. Objectives and methodology

Setting. We consider a multivariate regression problem of observations $y \in \mathbb{R}^m$, where $m \ge 3$, based on features $x \in$ \mathbb{R}^d , where $d \ge 1$. The observations enjoy some hierarchical structure: some of their components (henceforth referred 105 to as aggregated levels) are given by sums over subsets 106 of other components (henceforth referred to as the most disaggregated levels). More formally, up to reordering the components of y, there exist $2 \leq n < m$ and a $m \times n$ 109



Figure 1. Two examples of hierarchical structures and their associated structural matrices H.

matrix H of the form

$$H = \begin{bmatrix} \mathrm{Id}_n \\ H_{\mathrm{sub}} \end{bmatrix} \quad \text{such that} \quad \boldsymbol{y} = H \boldsymbol{y}_{1:n} \,, \qquad (1)$$

where H_{sub} is any $(m-n) \times n$ matrix of real numbers. The matrix H encoding the hierarchical summation constraints is called the structural¹ matrix. Two examples are provided in Figure 1, of the "bottom-up" form (i.e., with matrices $H_{\rm sub}$ of some specific form, but we recall that we will require no specific assumption on H_{sub}).

Definition 2.1. Vectors $u \in \mathbb{R}^m$ satisfying the linear constraints $u = Hu_{1:n}$ are called coherent. The subspace Im(H) of all such vectors is called the coherent subspace.

2.1. Objectives

We assume that i.i.d. data $(\boldsymbol{x}_t, \boldsymbol{y}_t)_{1 \leq t \leq T}$ is available to perform the regression task (we do so for the sake of exposition; Assumption 3.1 will later relax this requirement). We now describe the objective(s) targeted: constructing prediction rectangles with a coverage objective (\star) and a small-length objective $(\star\star)$.

The primary objective is to construct prediction sets based on this T-sample, of the form $C_1 \times \ldots \times C_m$, where each $C_i : \boldsymbol{x} \in \mathbb{R}^d \mapsto C_i(\boldsymbol{x})$ is an application taking subsets of \mathbb{R} as values. Consider a new data point $(\boldsymbol{x}_{T+1}, \boldsymbol{y}_{T+1})$ i.i.d. from the T-sample. The prediction sets should be such that each component $y_{T+1,i}$ of the observations \boldsymbol{y}_{T+1} be predicted by $C_i(\boldsymbol{x}_{T+1})$ with a coverage level of approximatively $1 - \alpha$:

$$\forall i \in [m], \quad \mathbb{P}\big(y_{T+1,i} \in C_i(\boldsymbol{x}_{T+1})\big) \approx 1 - \alpha, \quad (\star)$$

where the probability \mathbb{P} is with respect to both $(\boldsymbol{x}_{T+1}, \boldsymbol{y}_{T+1})$ and $(\boldsymbol{x}_t, \boldsymbol{y}_t)_{1 \leqslant t \leqslant T}$.

¹In the literature of forecast reconciliation, this matrix is usually denoted by S. We rather keep this letter for non-conformity scores. Also, the most disaggregated levels are often the last n components, while we consider the first n components.

110 This objective is different from the typical objective in other 111 contributions on multivariate conformal prediction (see Mes-112 soudi et al., 2021, Messoudi et al., 2022, Feldman et al., 113 2023), which is about a global (not component-wise) cov-114 erage guarantee: output a prediction region $C(x_{T+1})$ such 115 that

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143 144 145

146

147

148

149 150

151

$$\mathbb{P}(\boldsymbol{y}_{T+1} \in C(\boldsymbol{x}_{T+1})) \approx 1 - \alpha.$$
(2)

We are interested in the objective (\star) because we want to possibly concentrate on specific nodes within the hierarchy and get individual coverage guarantees for them. For instance, in our real-data application (see Section 4), this could correspond to making sure that each recharging station is well-dimensioned.

In addition, we discuss a simple way of leveraging the approach developed later in this article for the joint-coverage objective (2): see Appendix E.

A secondary objective is to ensure that the prediction sets output are efficient, i.e., are as small as possible. We will be mostly interested in rectangles $C_1(\mathbf{x}) \times \ldots \times C_m(\mathbf{x})$ based on intervals $C_i(\mathbf{x}) = [\mu_i(\mathbf{x}) + a_i(\mathbf{x}), \mu_i(\mathbf{x}) + b_i(\mathbf{x})]$, whose respective lengths are denoted by

$$\ell(C_i(\boldsymbol{x})) = b_i(\boldsymbol{x}) - a_i(\boldsymbol{x})$$

One quantification of the size of such a rectangle is given by

$$\sum_{i=1}^m w_i \, \ell \big(C_i(\boldsymbol{x}) \big)^2 \,,$$

where we fix some vector $\boldsymbol{w} = (w_1, \dots, w_m)^{\mathsf{T}}$ of positive numbers. The weights \boldsymbol{w} may be used to ponder the components based on their respective importance. The secondary objective then formally corresponds to

minimizing
$$\mathbb{E}\left[\sum_{i=1}^{m} w_i \ell (C_i(\boldsymbol{x}_{T+1}))^2\right], \quad (\star\star)$$

where, again, the expectation \mathbb{E} is with respect to both $(\boldsymbol{x}_{T+1}, \boldsymbol{y}_{T+1})$ and $(\boldsymbol{x}_t, \boldsymbol{y}_t)_{1 \leq t \leq T}$.

2.2. Hierarchical Split Conformal Prediction

Positioning. Lei et al. (2018) introduced a variant of con-152 formal prediction called split conformal prediction [SCP], 153 a procedure based on splitting data between a training set 154 155 (to learn a regressor function) and a calibration set (to compute estimation errors, a.k.a. residuals or non-conformity 156 scores). This procedure has been extensively studied in the 157 univariate case, and often through considering the absolute 158 values of the non-conformity scores, which leads to centered 159 160 intervals. We are interested in two extensions of this basic setting: multivariate SCP and signed non-conformity scores. 161

Signed non-conformity scores were already considered in the univariate case by Linusson et al. (2014). They are handy

in our setting because we consider linear constraints: the signed non-conformity scores $\hat{s} = y - \hat{y}$ between coherent observations y and forecasts \hat{y} are also coherent, while the vector of their absolute values is not coherent in general.

Multivariate SCP was already studied by Messoudi et al. (2021) but with somewhat different objectives: in particular, the design of prediction regions for the vectors y_{T+1} , while our objective (**) is about separate prediction intervals for each component. A more-in-depth discussion of the links and differences between the objectives of Messoudi et al. (2021) and (**) may be found in Appendix E.

Formal description: plain multivariate version. Formally, the component-wise coverage objectives (\star) may be achieved by a component-wise SCP, as follows. Data splitting corresponds to partitioning [T] into the subsets $\mathcal{D}_{\text{train}}$ and \mathcal{D}_{calib} , with respective cardinalities denoted by T_{train} and T_{calib} . With pairs $(\boldsymbol{x}_t, \boldsymbol{y}_t)$ indexed by $t \in \mathcal{D}_{\text{train}}$, a regressor function $\widehat{\mu} : x \in \mathbb{R}^d \mapsto \widehat{\mu}(x) \in \mathbb{R}^m$ is built, thanks to some regression algorithm \mathcal{A} provided as input parameter to the SCP procedure. Then, on the calibration set, i.e., for each $t \in \mathcal{D}_{\text{calib}}$, point estimates $\widehat{y}_t = \widehat{\mu}(x_t)$ and signed non-conformity scores (also known as signed estimation errors or signed residuals) $\widehat{s}_t = y_t - \widehat{y}_t$ are computed. The component-wise character of the procedure appears in the third and final step, where prediction intervals are output component by component. Indeed, for each component $i \in [m]$, we order separately the *i*-th components $\hat{s}_{t,i}$ of the non-conformity scores \hat{s}_t , where $t \in \mathcal{D}_{calib}$; the ordered values are denoted as follows, by using the standard notation of order statistics:

$$\widehat{s}_{(1),i} \leqslant \ldots \leqslant \widehat{s}_{(T_{\text{calib}}),i}$$
.

We also define $\hat{s}_{(0),i} = -\infty$ and $\hat{s}_{(T_{\text{calib}}+1),i} = +\infty$. Finally, the prediction interval for $y_{T+1,i}$ based on the corresponding features x_{T+1} is

$$\widehat{C}_{i}(\boldsymbol{x}_{T+1}) = \left[\widehat{\mu}_{i}(\boldsymbol{x}_{T+1}) + \widehat{s}_{\left(\lfloor (T_{\text{calib}}+1)\alpha/2 \rfloor\right), i}, \\ \widehat{\mu}_{i}(\boldsymbol{x}_{T+1}) + \widehat{s}_{\left(\lceil (T_{\text{calib}}+1)(1-\alpha/2) \rceil\right), i}\right],$$

where $1 - \alpha \in (0, 1)$ is the confidence level targeted and where $\hat{\mu}_i(\boldsymbol{x}_{T+1})$ is the *i*-th component of the point estimate $\hat{\boldsymbol{\mu}}(\boldsymbol{x}_{T+1})$. The prediction rectangle for \boldsymbol{y}_{T+1} is the Cartesian product of the prediction intervals $\hat{C}_i(\boldsymbol{x}_{T+1})$.

The thus-defined plain multivariate version of SCP with signed non-conformity scores is summarized in Algorithm 1. We use it as a benchmark and now introduce a generalization of this algorithm taking the hierarchical structure H into account.

Alg	orithm 1 Plain multivariate SCP with signed scores
Pai	rameters: confidence level $1 - \alpha$; regression algo-
	rithm \mathcal{A} ; partition of $[T]$ into subsets \mathcal{D}_{train} and \mathcal{D}_{calib}
	of respective cardinalities T_{train} and T_{calib}
1:	Build the regressor $\widehat{\mu}(\cdot) = \mathcal{A}ig((m{x}_t,m{y}_t)_{t\in\mathcal{D}_{ ext{train}}}ig)$
2:	Denote $\widehat{\boldsymbol{\mu}}(\cdot) = (\widehat{\mu}_1(\cdot), \dots, \widehat{\mu}_m(\cdot))$
3:	for $t \in \mathcal{D}_{\text{calib}}$ let $\widehat{m{y}}_t = \widehat{m{\mu}}(m{x}_t)$ and $\widehat{m{s}}_t = m{y}_t - \widehat{m{y}}_t$
4:	for each $i \in [m]$ do
5:	order the $(\hat{s}_{t,i})_{t \in \mathcal{D}_{\text{calib}}}$ into $\hat{s}_{(1),i} \leq \ldots \leq \hat{s}_{(T_{\text{calib}}),i}$
	and define $\widehat{s}_{(0),i} = -\infty$ and $\widehat{s}_{(T_{\text{calib}}+1),i} = +\infty$
6:	$\operatorname{let} \widehat{q}_{\alpha/2}^{(i)} = \widehat{s}_{\left(\lfloor (T_{\operatorname{calib}} + 1)\alpha/2 \rfloor \right), i}$
	and $\widehat{q}_{1-\alpha/2}^{(i)} = \widehat{s}_{\left(\left\lceil (T_{\text{calib}}+1)(1-\alpha/2) \right\rceil \right),i}$
7:	$\operatorname{set} \widehat{C}_i(\cdot) = \left[\widehat{\mu}_i(\cdot) + \widehat{q}_{\alpha/2}^{(i)}, \widehat{\mu}_i(\cdot) + \widehat{q}_{1-\alpha/2}^{(i)}\right]$
8:	return $\widehat{C}_1(\boldsymbol{x}_{T+1}), \ldots, \widehat{C}_m(\boldsymbol{x}_{T+1})$
Par	cameters: confidence level $1 - \alpha$; regression algo-
	of respective cardinalities T_{train} and T_{calib} ; matrix P
1:	Build the regressor $\widehat{\mu}(\cdot) = \mathcal{A}((\boldsymbol{x}_t, \boldsymbol{y}_t)_{t \in \mathcal{D}_{\text{train}}})$
2:	Denote $P\widehat{\mu}(\cdot) = \widetilde{\mu}(\cdot) = (\widetilde{\mu}_1(\cdot), \dots, \widetilde{\mu}_m(\cdot))$
3:	for $t \in \mathcal{D}_{ ext{calib}}$ let $\widetilde{m{y}}_t = \widetilde{m{\mu}}(m{x}_t)$ and $\widetilde{m{s}}_t = m{y}_t - \widetilde{m{y}}_t$
4:	for each $i \in [m]$ do
5:	order the $(\tilde{s}_{t,i})_{t \in \mathcal{D}_{calib}}$ into $\tilde{s}_{(1),i} \leq \ldots \leq \tilde{s}_{(T_{calib}),i}$
	and define $\widetilde{s}_{(0),i} = -\infty$ and $\widetilde{s}_{(T_{\text{calib}}+1),i} = +\infty$
6:	$\operatorname{let} \widetilde{q}_{\alpha/2}^{(i)} = \widetilde{s}_{\left(\lfloor (T_{\operatorname{calib}}+1)\alpha/2 \rfloor\right), i}$
	and $\widetilde{q}_{1-\alpha/2}^{(i)} = \widetilde{s}_{\left[\left[(T_{\text{calib}}+1)(1-\alpha/2)\right]\right],i}$
7:	set $\widetilde{C}_i(\cdot) = \left[\widetilde{\mu}_i(\cdot) + \widetilde{q}_{\alpha/2}^{(i)}, \widetilde{\mu}_i(\cdot) + \widetilde{q}_{1-\alpha/2}^{(i)}\right]$
8:	$\widetilde{C}(m) = \widetilde{C}(m)$
0.	return $C_1(x_{T+1}), \ldots, C_m(x_{T+1})$

Formal description: hierarchical version of SCP. The hierarchical version of SCP is stated in Algorithm 2 and only differs from the plain multivariate version stated as Algorithm 1 in line 2, where a projection matrix P onto the coherent subspace Im(H) should be used: the regressor function considered is $\tilde{\mu} = P\hat{\mu}$, instead of simply $\hat{\mu}$, and thus outputs point estimates that are coherent in the case where $\text{Im}(P) \subseteq \text{Im}(H)$. The rest of the procedure is similar.

Algorithm 1 is a special case of Algorithm 2, for the choice $P = \text{Id}_m$. We however provide two separate statements to clarify the notation: $\widehat{}$ -type quantities are for the plain multivariate version (Algorithm 1), which we use as a benchmark, while $\widetilde{}$ -type quantities refer to the hierarchical version (Algorithm 2) to be used with a projection matrix Ponto Im(H). **Projection matrices.** As indicated, we will be mostly interested in projection matrices P onto Im(H) for Algorithm 2. We will most often take them of the form $P_W = H(H^{\top}WH)^{-1}H^{\top}W$, where W is a symmetric definite positive matrix; see Lemma B.5 in Appendix B for a proof that the P_W are indeed projections onto Im(H), as well as a statement of additional properties that they enjoy.

3. Statements of the theoretical results

In conformal prediction, results typically hold in great generality. In particular, we will require no direct assumption on the regression algorithm \mathcal{A} , which will be treated as a black-box regression procedure that does not even have to output coherent point estimates (hence the consideration of a projection matrix P in Algorithm 2). No assumption will be required on the data split between $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{calib}}$, but the following key assumption of i.i.d. behavior will be issued on the data (or only on the non-conformity scores).

Assumption 3.1. The non-conformity scores $\hat{s}_t = y_t - \hat{y}_t$, for $t \in \mathcal{D}_{\text{calib}} \cup \{T+1\}$ of the plain multivariate version of SCP (Algorithm 1) are i.i.d. This is in particular the case when data $(x_t, y_t)_{1 \leq t \leq T+1}$ is i.i.d.

The second part of Assumption 3.1 follows from the fact that $\hat{s}_t = y_t - \hat{\mu}(x_t)$, where $\hat{\mu}$ only depends on the data in $\mathcal{D}_{\text{train}}$ and is therefore independent from the data (x_t, y_t) for $t \in \mathcal{D}_{\text{calib}} \cup \{T+1\}$.

3.1. Coverage objective (*)

The theorem below is related to objective (*). It of course holds for Algorithm 1, given that it is a special case of Algorithm 2 with $P = \text{Id}_m$. At this stage, no assumption is required on the matrix P. The standard proof (together with references to earlier similar proofs) may be found in Appendix A.

Theorem 3.2 (Coverage). Fix $\alpha \in (0, 1)$. Algorithm 2, used with any regression algorithm A and any matrix P such that PH = H, ensures that whenever Assumption 3.1 (i.i.d. scores) holds,

$$\forall i \in [m], \quad \mathbb{P}(y_{T+1,i} \in \tilde{C}_i(\boldsymbol{x}_{T+1})) \ge 1 - \alpha.$$

In addition, if the non-conformity scores $(\hat{s}_t)_{t \in \mathcal{D}_{\text{calib}} \cup \{T+1\}}$ are almost surely distinct, then

$$\forall i \in [m], \quad \mathbb{P}\left(y_{T+1,i} \in \widetilde{C}_i(\boldsymbol{x}_{T+1})\right) \leq 1 - \alpha + \frac{2}{T_{\text{calib}} + 1}$$

3.2. Small-length objective (******), weak version

We now move on to the small-length objective $(\star\star)$, and start with a weak version thereof, relying on one fixed weight vector w (see next section for a stronger version).

220 As mentioned above, we issue no direct assumption on 221 the regression algorithm \mathcal{A} . However, we output a mild 222 assumption on the distribution of the non-conformity scores, 223 which may be seen as un indirect assumption on A. This 224 assumption falls within the model-agnostic gist of conformal 225 prediction and is standard in the literature of probabilistic forecast reconciliation (see, e.g., Panagiotelis et al., 2023 227 or Wickramasuriya, 2024). It states that non-conformity 228 scores follow a distribution, called an elliptical distribution, 229 derived from a spherical distribution; for more details, see 230 Appendix B (which in turns refers to Kollo & von Rosen, 231 2005, Section 2.3).

Definition 3.3. A random vector z follows a spherical distribution over \mathbb{R}^k if z and Γz have the same distribution for all $k \times k$ orthogonal matrices Γ .

232

233

234

240

241

242

243

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

269

270 271

272

273

274

235 **Definition 3.4.** An elliptical distribution over \mathbb{R}^m is of the 236 form c + Mz, for a deterministic vector $c \in \mathbb{R}^m$, a $m \times k$ 238 matrix M such that MM^{\top} has rank k, and a random vector 239 z following a spherical distribution over \mathbb{R}^k .

A given spherical distribution thus generates a family \mathcal{F} of elliptical distributions enjoying a stability property through linear transformations.

Example 3.5. The simplest example of elliptical distributions consists of multivariate normal distributions (which are light-tailed distributions). Other examples include multivariate *t*-distributions and symmetric multivariate Laplace distributions (both heavy tailed) and the uniform distribution on an ellipse (no tail).

We are now ready to state the key assumption used to establish in Theorem 3.7 that the hierarchical version of SCP performs better than its plain multivariate version.

Assumption 3.6. The (i.i.d.) non-conformity scores $\hat{s}_t = y_t - \hat{y}_t$, for $t \in \mathcal{D}_{calib}$, of the plain multivariate version of SCP (Algorithm 1) follow some elliptical distribution (whose parameters are unknown to the learner).

Theorem 3.7. Let $w \in \mathbb{R}^m$ be a vector of positive numbers. Under Assumptions 3.1 and 3.6 (i.i.d. scores with elliptical distribution), the hierarchical version of SCP (Algorithm 2) run with $P = P_w$, where

$$P_{\boldsymbol{w}} \stackrel{\text{def}}{=} H \left(H^{\top} \operatorname{diag}(\boldsymbol{w}) H \right)^{-1} H^{\top} \operatorname{diag}(\boldsymbol{w}) , \qquad (3)$$

provides prediction rectangles that are more efficient than the ones output by the plain multivariate version of SCP (Algorithm 1) in the following sense:

$$\mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell\big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big)^2\right] \leqslant \mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell\big(\widehat{C}_i(\boldsymbol{x}_{T+1})\big)^2\right].$$

Sketch of proof. The centered scores $\tilde{s}_t - \mathbb{E}[\tilde{s}_t]$ are i.i.d. according to a centered elliptical distribution as $t \in \mathcal{D}_{\text{calib}}$.

Thus, their *i*-th components have the same distribution ν up to a scaling factor denoted by $\sqrt{\gamma_i}$. Let $(v_t)_{t \in \mathcal{D}_{calib}}$ be i.i.d. variables distributed according to ν , consider their order statistics $v_{(1)} \leq \ldots \leq v_{(T_{calib})}$, and set

$$L_{\alpha} = v_{\left(\lceil (T_{\text{calib}}+1)(1-\alpha/2) \rceil \right)} - v_{\left(\lfloor (T_{\text{calib}}+1)\alpha/2 \rfloor \right)}.$$

We thus have that $\ell(C_i(\boldsymbol{x}_{T+1}))$ has the same distribution as $\sqrt{\gamma_i} L_{\alpha}$, for each $i \in [m]$. Therefore,

$$\mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell\big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big)^2\right] = \mathbb{E}\left[L_{\alpha}^2\right] \sum_{i \in [m]} w_i \gamma_i \,. \tag{4}$$

It may be shown that γ_i is the (i, i)-th element of a matrix of the form $P_{\boldsymbol{w}} \Gamma P_{\boldsymbol{w}}^{\top}$, where Γ is symmetric positive semidefinite. A similar result holds for the \hat{C}_i , with scaling factors given by the diagonal elements of Γ . It thus suffices to show that

$$\sum_{i \in [m]} w_i \Gamma_{i,i} \geqslant \sum_{i \in [m]} w_i \left(P_{\boldsymbol{w}} \Gamma P_{\boldsymbol{w}}^{\mathsf{T}} \right)_{i,i}, \qquad (5)$$

i.e., that $\operatorname{Tr}(\operatorname{diag}(w)\Gamma) \geq \operatorname{Tr}(\operatorname{diag}(w)P_w\Gamma P_w^{\top})$. The latter result essentially follows from taking expectations of a Pythagorean inequality, and is a result of our own, though inspired by the literature of forecast reconciliation.

The complete proof may be found in Appendix B, including a justification of why P_w is well defined.

Some words on the challenges overcome. We sketched above how we connected the tools of conformal prediction and of forecast reconciliation.

The difficulty mostly lied in the component-wise approach, imposed by the component-wise coverage guarantees (\star) targeted. Indeed, Appendix E.2 explains in detail why efficiency results are straightforward when some joint coverage is targeted.

More precisely (and as detailed in Appendix D.2), the main blocking point was to relate the minimization of squared lengths (4) to what may be rephrased as some trace minimization (5). Such relationships are classic in the literature of forecast reconciliation, but they rely on assumptions of unbiasedness (i.e., of centered non-conformity scores, which we are not ready to consider). However, by resorting to signed scores, the expectations of the scores cancel out: the distribution of L_{α} is stable when the $(v_t)_{t \in \mathcal{D}_{calib}}$ are translated.

3.3. Objective $(\star\star)$: stronger but oracle version

We resort to tools from the theory of forecast reconciliation to improve the results of Theorem 3.7 and have them hold simultaneously for all possible positive weight vectors w. 275 However, this improvement is only for an oracle strategy 276 relying on a projection matrix $P_{\Sigma^{-1}}$ depending on the co-277 variance matrix Σ of the scores (unknown to the learner). 278 This is why we state next (see Algorithm 3) a "practical"

implementation of this oracle, which adds an estimation step for Σ to the hierarchical SCP procedure stated as Algorithm 2.

Assumption 3.8. The (i.i.d.) non-conformity scores $\hat{s}_t = y_t - \hat{y}_t$, for $t \in \mathcal{D}_{\text{calib}}$, of the plain multivariate version of SCP (Algorithm 1) have a bounded second-order moment. We denote by Σ their (positive definite) covariance matrix.

The projection $P_{\Sigma^{-1}}$ to be considered in Theorem 3.10 is the one used in the Minimum-Trace projection in forecast reconciliation (Wickramasuriya et al., 2019). The name comes from the following optimality result, for which we provide an elementary proof in Appendix C (this proof may be considered in itself as a result of interest). Let

295

296

297

299

300 301

302

303

304

329

$$P_{\Sigma^{-1}} \stackrel{\text{def}}{=} H \left(H^{\top} \Sigma^{-1} H \right)^{-1} H^{\top} \Sigma^{-1} .$$
 (6)

Lemma 3.9 (Minimum-Trace projection). Let W and Σ be two symmetric $m \times m$ matrices, where W positive semidefinite and Σ is positive definite. Then, for all projection matrices P onto Im(H),

$$\operatorname{Tr}(WP_{\Sigma^{-1}}\Sigma P_{\Sigma^{-1}}^{\top}) \leq \operatorname{Tr}(WP\Sigma P^{\top}).$$

We now state our main results, Theorem 3.10 and Corollary 3.11. Their proofs may be found in Appendix C and consist of direct adaptations of the proof of Theorem 3.7, together with an application of Lemma 3.9.

To do so, we denote by $\tilde{C}_1^{\star}(\boldsymbol{x}_{T+1}) \times \ldots \times \tilde{C}_m^{\star}(\boldsymbol{x}_{T+1})$ the prediction rectangles output by the hierarchical version of SCP (Algorithm 2) run with $P = P_{\Sigma^{-1}}$ defined in (6), and keep the notation $\tilde{C}_1(\boldsymbol{x}_{T+1}) \times \ldots \times \tilde{C}_m(\boldsymbol{x}_{T+1})$ for the prediction rectangles output by the same strategy run with any other choice of a projection matrix.

We obtain the following efficiency result, which is actually
stronger than the objective (**) stated, as the comparison
holds component-wise.

Theorem 3.10. Under Assumptions 3.1, 3.6, and 3.8 (i.i.d. scores with elliptical distribution admitting a second-order moment), the hierarchical version of SCP (Algorithm 2) run with $P = P_{\Sigma^{-1}}$ provides prediction rectangles more efficient than with any other choice of a projection matrix onto Im(H):

327
$$\forall i \in [m], \quad \mathbb{E}\left[\ell\left(\widetilde{C}_{i}^{\star}(\boldsymbol{x}_{T+1})\right)^{2}\right] \leq \mathbb{E}\left[\ell\left(\widetilde{C}_{i}(\boldsymbol{x}_{T+1})\right)^{2}\right].$$
328

Corollary 3.11. *In the setting and under the assumptions of Theorem 3.10, we also have*

$$\forall i \in [m], \quad \mathbb{E}\left[\ell\left(\widetilde{C}_{i}^{\star}(\boldsymbol{x}_{T+1})\right)^{2}\right] \leq \mathbb{E}\left[\ell\left(\widehat{C}_{i}(\boldsymbol{x}_{T+1})\right)^{2}\right],$$

where the prediction rectangle $\widehat{C}_1(\boldsymbol{x}_{T+1}) \times \ldots \times \widehat{C}_m(\boldsymbol{x}_{T+1})$ is output by the plain multivariate version of SCP (Algorithm 1).

The component-wise comparisons stated above in Theorem 3.10 and Corollary 3.11 correspond to inequalities similar to the ones of Theorem 3.7 holding for all positive weight vectors w:

$$\forall \boldsymbol{w} \in (0, +\infty)^m, \quad \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell\big(\widetilde{C}_i^{\star}(\boldsymbol{x}_{T+1})\big)^2\right]$$
$$\leqslant \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell\big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big)^2\right].$$

(and not just for a single vector w as in Theorem 3.7).

However, the covariance matrix Σ is unknown and therefore, running Algorithm 2 with $P_{\Sigma^{-1}}$ is an oracle strategy. We turn it into a practical strategy by adding an estimation step.

3.4. Hierarchical SCP with estimated covariance matrix

Theorem 3.10 advocates using Algorithm 2 with $P_{\Sigma^{-1}}$ but the covariance matrix Σ of the (unprojected) non-conformity scores is unknown. A natural idea is of course to estimate it: this is why data is now split into three subsets $\mathcal{D}_{\text{train}}$ (train set), $\mathcal{D}_{\text{estim}}$ (estimation set), and $\mathcal{D}_{\text{calib}}$ (calibration set) of respective cardinalities T_{train} , T_{estim} and T_{calib} . A regression function $\hat{\mu}$ is still learned on data of $\mathcal{D}_{\text{train}}$, in some black-box way. The covariance matrix of the non-conformity scores built based on $\hat{\mu}$ is estimated on a fraction of the remaining data, indexed by $\mathcal{D}_{\text{estim}}$. This estimate then determines a projection matrix P to be used to run the final part of Algorithm 2: project scores, rank their components, and deduce the prediction intervals of each component based on the latter rankings.

More precisely, based on the scores $\hat{s}_t = y_t - \hat{\mu}(x_t)$, for $t \in \mathcal{D}_{\text{estim}}$, we compute the sample mean \bar{s} and sample covariance matrix $\hat{\Sigma}$:

$$\overline{s} = \frac{1}{T_{\text{estim}}} \sum_{t \in \mathcal{D}_{\text{estim}}} \widehat{s}_t \,, \quad \widehat{\Sigma} = \frac{1}{T_{\text{estim}}} \sum_{t \in \mathcal{D}_{\text{estim}}} \left(\widehat{s}_t - \overline{s} \right) \left(\widehat{s}_t - \overline{s} \right)^\top.$$

We assume that $\widehat{\Sigma}$ is symmetric positive definite, which was always the case in our simulations. A projection matrix $P = \mathcal{P}(\widehat{\Sigma})$ onto $\operatorname{Im}(H)$ is then computed based on $\widehat{\Sigma}$, through some function \mathcal{P} that operates on symmetric positive definite $m \times m$ matrices. Examples of such functions include

$$\mathcal{P}_{\mathrm{MinT}}: M \longmapsto H \left(H^{\mathsf{T}} M^{-1} H \right)^{-1} H^{\mathsf{T}} M^{-1} , \qquad (7)$$

330 Algorithm 3 Hierarchical SCP with estimated covariances **Parameters:** confidence level $1 - \alpha$; regression algorithm \mathcal{A} ; partition of [T] into three subsets $\mathcal{D}_{\text{train}}$, $\mathcal{D}_{\text{estim}}$, 333 and \mathcal{D}_{calib} of respective cardinalities T_{train} , T_{estim} and 334 T_{calib} ; function \mathcal{P} operating on $m \times m$ matrices 335 1: Build the regressor $\widehat{\mu}(\cdot) = \mathcal{A}((\boldsymbol{x}_t, \boldsymbol{y}_t)_{t \in \mathcal{D}_{\text{train}}})$ {----Modifications w.r.t. Algorithm 2 start here-----} 337 2: for $t \in \mathcal{D}_{\text{estim}}$ let $\widehat{y}_t = \widehat{\mu}(x_t)$ and $\widehat{s}_t = y_t - \widehat{y}_t$ 338 3: Compute the sample mean of the $(\hat{s}_t)_{t \in \mathcal{D}_{estim}}$, 339 $\overline{oldsymbol{s}} = rac{1}{T_{ ext{estim}}} \sum_{t \in \mathcal{D}_{ ext{estim}}} \widehat{oldsymbol{s}}_t,$ 340 and their 341 342 covariance matrix $\widehat{\Sigma} = \frac{1}{T_{\text{estim}}} \sum_{t \in \mathcal{D}} (\widehat{s}_t - \overline{s}) (\widehat{s}_t - \overline{s})^{\top}$ 343 344 4: Let $P = \mathcal{P}(\widehat{\Sigma})$ 345 {----Modifications w.r.t. Algorithm 2 stop here-----} 346 347 5: Denote $P\widehat{\boldsymbol{\mu}}(\cdot) = \widetilde{\boldsymbol{\mu}}(\cdot) = (\widetilde{\mu}_1(\cdot), \dots, \widetilde{\mu}_m(\cdot))$ 6: for $t \in \mathcal{D}_{calib}$ let $\widetilde{\boldsymbol{y}}_t = \widetilde{\boldsymbol{\mu}}(\boldsymbol{x}_t)$ and $\widetilde{\boldsymbol{s}}_t = \boldsymbol{y}_t - \widetilde{\boldsymbol{y}}_t$ 349 7: for each $i \in [m]$ do 350 order the $(\tilde{s}_{t,i})_{t \in \mathcal{D}_{\text{calib}}}$ into $\tilde{s}_{(1),i} \leq \ldots \leq \tilde{s}_{(T_{\text{calib}}),i}$ and define $\tilde{s}_{(0),i} = -\infty$ and $\tilde{s}_{(T_{\text{calib}}+1),i} = +\infty$ 8: 351 352 $\det \widetilde{q}_{\alpha/2}^{(i)} = \widetilde{s}_{\left(\lfloor (T_{\mathrm{calib}}+1)\alpha/2 \rfloor \right), i}$ 9: 353 and $\widetilde{q}_{1-\alpha/2}^{(i)} = \widetilde{s}_{\left(\left\lceil (T_{\text{calib}}+1)(1-\alpha/2) \right\rceil \right),i}$ 354 355 set $\widetilde{C}_i(\cdot) = \left[\widetilde{\mu}_i(\cdot) + \widetilde{q}_{\alpha/2}^{(i)}, \widetilde{\mu}_i(\cdot) + \widetilde{q}_{1-\alpha/2}^{(i)}\right]$ 356 10: 357 11: return $\widetilde{C}_1(\boldsymbol{x}_{T+1}), \ldots, \widetilde{C}_m(\boldsymbol{x}_{T+1})$ 358 359 360

which mimics the expression for $P_{\Sigma^{-1}}$ in Theorem 3.10, corresponding to the Min-Trace projection, hence the MinT subscript. Other examples are discussed below.

361

362

363

365

366

367

368

380

381

382

383

384

The procedure is summarized in Algorithm 3. The statement of the latter only differs from the one of Algorithm 2 through the additional lines 2–3–4, which consider the estimation set \mathcal{D}_{estim} and build the projection P based on $\hat{\Sigma}$.

369 Coverage guarantees. The coverage guarantees of 370 Theorem 3.2 hold also for Algorithm 3 when data 371 $({m x}_t, {m y}_t)_{1\leqslant t\leqslant T+1}$ is i.i.d. Indeed, the non-conformity scores 372 $\widetilde{s}_t = y_t - \widetilde{y}_t$, for $t \in \mathcal{D}_{calib} \cup \{T+1\}$, are then still i.i.d., 373 which is the only property needed in the proof of Theo-374 rem 3.2 (located in Appendix A; see the end of its first 375 paragraph). This follows from the fact that these scores are defined based on $\hat{\mu}$, on $\hat{\Sigma}$, and on $(\boldsymbol{x}_t, \boldsymbol{y}_t)_{t \in \mathcal{D}_{\text{calib}} \cup \{T+1\}}$, 376 377 where the latter are independent from $\hat{\mu}$ and $\hat{\Sigma}$, which only 378 depend on data of \mathcal{D}_{train} and \mathcal{D}_{estim} . 379

Other examples of projection functions \mathcal{P} . When data is scarce, the estimates $\hat{\Sigma}$ may be poor, which would cause issues when taking its inverse as in (7). A more robust approach is to consider the vector of the inverses of the

Table 1. Summary of the algorithms implemented: nicknames and corresponding formal definitions, through the algorithm number, the required parameter P or \mathcal{P} (if applicable), and the defining equation of the latter. The first five algorithms are practical algorithms, while the sixth algorithm is an oracle.

Nickname	Algorithm	Parameter	Equation
Direct	Alg. 1		
OLS	Alg. 2	P_1	Eq. $(3) + (9)$
MinT	Alg. 3	$\mathcal{P}_{\mathrm{MinT}}$	Eq. (7)
WLS	Alg. 3	$\mathcal{P}_{\mathrm{WLS}}$	Eq. (8)
Combi	Alg. 3	$\mathcal{P}_{\mathrm{Combi}}$	Eq. (10)
Oracle MinT	Alg. 2	$P_{\Sigma^{-1}}$	Eq. (6)

diagonal elements of $\hat{\Sigma}$ only (i.e., the vector of the inverses of the variances of the components of the scores) and the associated projection matrix P_w , as in Theorem 3.7. This corresponds to some data-based weighted least squares [WLS], as pointed out by Hyndman et al. (2016):

$$\mathcal{P}_{\mathrm{WLS}}: M \longmapsto H \left(H^{\mathsf{T}} D_M^{-1} H \right)^{-1} H^{\mathsf{T}} D_M^{-1}, \qquad (8)$$

where $D_M = \text{diag}((M_{i,i})_{i \in [m]})$. In contrast, an ordinary least squares [OLS] approach (considered, for instance, in Hyndman et al., 2011) would be based on the orthogonal projection onto Im(H), which corresponds to P_1 in (3):

$$P_{\mathbf{1}} = H \left(H^{\top} H \right)^{-1} H^{\top}, \text{ where } \mathbf{1} = (1, \dots, 1)^{\top}.$$
 (9)

Another robust approach could be to use a combination (hence the subscript "Combi") of the \mathcal{P} functions defined above, as suggested by Hollyman et al. (2021), who argued that this combination does not have to be complicated to be efficient. We therefore consider

$$\mathcal{P}_{\text{Combi}}: M \longmapsto \frac{1}{3} \left(P_1 + \mathcal{P}_{\text{WLS}}(M) + \mathcal{P}_{\text{MinT}}(M) \right).$$
(10)

Lemma B.5 in Appendix B states that the \mathcal{P}_{MinT} , \mathcal{P}_{WLS} , and thus \mathcal{P}_{Combi} take projection matrices as values.

4. Simulations on artificial and real data

We compare the performance achieved by the algorithms presented in this article, as summarized in Table 1: we do so in terms of component-wise coverage levels and of lengths of the prediction intervals. We consider two simulation settings: one on artificially generated data, and one on real data consisting of daily energy demand to charge electric vehicles in Palo Alto, CA. In both cases, the hierarchical structure is the three-level hierarchy (with 8 nodes, of the form 5–2–1) described in the right-hand side of Figure 1.

Full details of the simulation settings (both for artificially generated data and real data) may be found in Appendix F. We only summarize below the main observations.



Figure 2. Component-wise coverage levels and prediction-interval lengths (top graphs) and total lengths (bottom graphs), for artificially generated data (*left graphs*) and real data of energy demand (*right graphs*). Empirical averages are reported, with standard errors (in both the x-axis and y-axis directions for top graphs). The layout of top graphs follows the tree-representation of the hierarchy.

Artificially generated data. We generate 1000 runs of 419 some experimental setting and obtain empirical estimates 420 of the component-wise coverage levels (x-axis positions in 421 the top graphs of Figure 2) and root-mean squared lengths 422 (y-axis positions of the same graphs). All algorithms obtain 423 component-wise coverage levels close to each other and 424 close to the target value $1 - \alpha = 90\%$. However, differ-425 ences are clearer in terms of component-wise lengths, with 426 our preferred strategy, MinT, i.e., Algorithm 3 with \mathcal{P}_{MinT} , 427 outperforming all others (as hoped from Theorem 3.10), 428 especially at the lower levels of the hierarchy. Simpler 429 strategies like OLS, i.e., Algorithm 2 with a plain Euclidean 430 projection, provide shorter lengths than Direct, the plain 431 SCP method, but are less efficient than MinT. These differ-432 ences in lengths are significant as far as the total lengths are 433 concerned: see the bottom graphs of Figure 2, where the 434 confidence intervals on the expected means are all disjoint. 435 The strategies are thus clearly ranked. 436

415

416

417 418

Real data of energy demand. The data consists of daily
 charging loads for 5 charging points in two different parking

lots (containing respectively 3 and 2 charging points, hence the stated 5–2–1 hierarchy). We consider that the daily data are i.i.d. and generate 360 runs of some experimental setting consisting, in particular, of selecting fractions of the data at random for the train, estimation, and calibration sets. We report empirical averages, as in the case of artificially generated data. All strategies obtain similar componentwise coverage guarantees, close to 90.5%. This outcome is in line with Theorem 3.2, which states that the coverage level is actually possibly slightly larger than $1-\alpha$, of a factor of order $1/T_{calib}$. The sample sizes are much smaller on this use case (T = 1780 vs. $T = 10^5$ for artificial data), which explains why the slightly larger coverage target appears clearly. Length-wise, the MinT strategy performs relatively poorly, which may be explained by the poor estimation of the covariance matrix (due to lack of data). The WLS and Combi strategies are more robust, as stated Section 3.4. The oracle version of MinT performs the best, as expected (but is an oracle).

440 Impact Statement

441

442

443

444

445

446

447

448

449 450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

472

473

474

475

476

481

482

483

484

485

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here. (The application to real-data of optimizing the charging of electric vehicles is an example of a positive societal consequence, but is only an example among many other possible ones.)

References

- Amara-Ouali, Y., Goude, Y., Massart, P., Poggi, J.-M., and Yan, H. A review of electric vehicle load open data and models. *Energies*, 14(8):2233, 2021.
- Amara-Ouali, Y., Goude, Y., Hamrouche, B., and Bishara, M. A benchmark of electric vehicle load and occupancy models for day-ahead forecasting on open charging session data. In *Proceedings of the Thirteenth ACM International Conference on Future Energy Systems (e-Energy'2022)*, pp. 193–207, 2022.
- Ando, S. and Narita, F. An alternative proof of minimum trace reconciliation. *Forecasting*, 6(2):456–461, 2024.
- Athanasopoulos, G., Ahmed, R. A., and Hyndman, R. J. Hierarchical forecasts for Australian domestic tourism. *International Journal of Forecasting*, 25(1):146–166, 2009.
- 469 Athanasopoulos, G., Hyndman, R. J., Kourentzes, N., and
 470 Panagiotelis, A. Forecast reconciliation: A review. *Inter-*471 *national Journal of Forecasting*, 40(2):430–456, 2024.
 - Duchi, J. C., Gupta, S., Jiang, K., and Sur, P. Predictive inference in multi-environment scenarios, 2024. Preprint, arXiv:2403.16336.
- 477 Dunn, R., Wasserman, L., and Ramdas, A. Distribution-free
 478 prediction sets for two-layer hierarchical models. *Journal*479 *of the American Statistical Association*, 118(544):2491–
 480 2502, 2023.
 - Feldman, S., Bates, S., and Romano, Y. Calibrated multipleoutput quantile regression with representation learning. *Journal of Machine Learning Research*, 24(24):1–48, 2023.
- Hollyman, R., Petropoulos, F., and Tipping, M. E. Understanding forecast reconciliation. *European Journal of Operational Research*, 294(1):149–160, 2021.
- Hyndman, R. J., Ahmed, R. A., Athanasopoulos, G., and
 Shang, H. L. Optimal combination forecasts for hierarchical time series. *Computational Statistics & Data Analysis*, 55(9):2579–2589, 2011.

- Hyndman, R. J., Lee, A. J., and Wang, E. Fast computation of reconciled forecasts for hierarchical and grouped time series. *Computational Statistics & Data Analysis*, 97: 16–32, 2016.
- Johnstone, C. and Cox, B. Conformal uncertainty sets for robust optimization. In *Proceedings of the Tenth Symposium on Conformal and Probabilistic Prediction with Applications (COPA'21)*, volume 152 of PMLR, pp. 72–90, 2021.
- Kollo, T. and von Rosen, D. Advanced Multivariate Statistics with Matrices. Mathematics and Its Applications. Springer, 2005.
- Lee, Y., Barber, R. F., and Willett, R. Distributionfree inference with hierarchical data, 2023. Preprint, arXiv:2306.06342.
- Lei, J., G'Sell, M., Rinaldo, A., Tibshirani, R. J., and Wasserman, L. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.
- Linusson, H., Johansson, U., and Löfström, T. Signed-error conformal regression. In Advances in Knowledge Discovery and Data Mining (PAKDD'2014), Part I, volume 8443 of Lecture Notes in Computer Science, pp. 224–236. Springer, 2014.
- Messoudi, S., Destercke, S., and Rousseau, S. Copula-based conformal prediction for multi-target regression. *Pattern Recognition*, 120:108101, 2021.
- Messoudi, S., Destercke, S., and Rousseau, S. Ellipsoidal conformal inference for multi-target regression. In *Proceedings of the Eleventh Symposium on Conformal and Probabilistic Prediction with Applications (COPA'22)*, volume 179 of PMLR, pp. 294–306, 2022.
- Panagiotelis, A., Athanasopoulos, G., Gamakumara, P., and Hyndman, R. J. Forecast reconciliation: A geometric view with new insights on bias correction. *International Journal of Forecasting*, 37(1):343–359, 2021.
- Panagiotelis, A., Gamakumara, P., Athanasopoulos, G., and Hyndman, R. J. Probabilistic forecast reconciliation: Properties, evaluation and score optimisation. *European Journal of Operational Research*, 306(2):693–706, 2023.
- Tibshirani, R. J., Foygel Barber, R., Candes, E., and Ramdas, A. Conformal prediction under covariate shift. In Advances in Neural Information Processing Systems (Neurips'2019), volume 32, pp. 2530–2540, 2019.
- Vovk, V., Gammerman, A., and Shafer, G. *Algorithmic Learning in a Random World*. Springer, 2005.

495 496 497	Wickramasuriya, S. L. Probabilistic forecast reconciliation under the Gaussian framework. <i>Journal of Business & Economic Statistics</i> , 42(1):272–285, 2024.	
498 499 500 501 502 503	Wickramasuriya, S. L., Athanasopoulos, G., and Hyndman, R. J. Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization. <i>Journal</i> of the American Statistical Association, 114(526):804– 819, 2019.	
504 505 506 507	Wood, S. <i>Package mgcv</i> , 2023. URL https://CRAN. R-project.org/package=mgcv. R package ver- sion 1.9-1.	
507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543	Wood, S. N. Generalized Additive Models: An Introduction with R. Chapman & Hall, 2017.	
544 545 546 547 548 549		

550 Appendices

The appendices contain the following material: first, detailed proofs of all claims made in the main body, namely,

- A proof of the coverage result (Theorem 3.2), in Appendix A;
- A proof of the efficiency result for fixed weights w (Theorem 3.7), in Appendix B;
- A proof of the component-wise efficiency results (Theorem 3.10 and Corollary 3.11), in Appendix C;

second, detailed reviews (literature review and formal description of classic results) of two topics:

- Forecast reconciliation, in Appendix D;
- Multi-target (multivariate) split conformal prediction, in Appendix E;

third,

• Full details for the simulations on artificial and real data, in Appendix F.

A. Proof of the coverage result (Theorem 3.2)

Theorem 3.2 (Coverage). Fix $\alpha \in (0, 1)$. Algorithm 2, used with any regression algorithm \mathcal{A} and any matrix P such that PH = H, ensures that whenever Assumption 3.1 (i.i.d. scores) holds,

$$\forall i \in [m], \quad \mathbb{P}(y_{T+1,i} \in \widetilde{C}_i(\boldsymbol{x}_{T+1})) \ge 1 - \alpha.$$

In addition, if the non-conformity scores $(\hat{s}_t)_{t \in \mathcal{D}_{calib} \cup \{T+1\}}$ are almost surely distinct, then

$$\forall i \in [m], \quad \mathbb{P} \big(y_{T+1,i} \in \widetilde{C}_i(\boldsymbol{x}_{T+1}) \big) \leqslant 1 - \alpha + \frac{2}{T_{\text{calib}} + 1}$$

The proof below uses a standard methodology in the literature of conformal prediction (see, for instance, Tibshirani et al., 2019, proof of Theorem 1), with rather immediate adaptations due to the multivariate context and to the choice of signed non-conformity scores.

Proof. The condition PH = H means that P leaves elements of Im(H) unchanged. Since observations y_t are coherent, we have, for all $t \in D_{calib}$,

$$\widetilde{\boldsymbol{s}}_t \stackrel{\text{def}}{=} \boldsymbol{y}_t - P \widehat{\boldsymbol{y}}_t = P (\boldsymbol{y}_t - \widehat{\boldsymbol{y}}_t) = P \widehat{\boldsymbol{s}}_t.$$

Assumption 3.1 thus entails that the non-conformity scores \tilde{s}_t , where $t \in \mathcal{D}_{calib} \cup \{T+1\}$, are also i.i.d., thus exchangeable—which is the only property we will use in the rest of this proof.

Fix $i \in [m]$. By definition of $\widetilde{C}_i(\boldsymbol{x}_{T+1})$ and of the score $\widetilde{\boldsymbol{s}}_{T+1} = \boldsymbol{y}_{T+1} - \widetilde{\boldsymbol{\mu}}(\boldsymbol{x}_{T+1})$, the event of interest may be rewritten as

$$\left\{y_{T+1,i} \in \widetilde{C}_i(\boldsymbol{x}_{T+1})\right\} = \left\{\widetilde{s}_{\left(\lfloor (T_{\text{calib}}+1)\alpha/2 \rfloor\right),i} \leqslant \widetilde{s}_{T+1,i} \leqslant \widetilde{s}_{\left(\lceil (T_{\text{calib}}+1)(1-\alpha/2)\rceil\right),i}\right\}.$$
(11)

If $\alpha \in (0,1)$ is so small that $(T_{\text{calib}}+1)\alpha/2 < 1$, i.e., $\alpha < 2/(T_{\text{calib}}+1)$, then $\widetilde{s}_{(\lfloor (T_{\text{calib}}+1)\alpha/2 \rfloor),i} = \widetilde{s}_{(0)} = -\infty$ and $\widetilde{s}_{(\lceil (T_{\text{calib}}+1)(1-\alpha/2)\rceil),i} = \widetilde{s}_{(T_{\text{calib}}+1)} = +\infty$. Thus,

$$\mathbb{P}(y_{T+1,i} \in C_i(\boldsymbol{x}_{T+1})) = 1$$

satisfies the claimed statements $\geq 1 - \alpha$ and $\leq 1 - \alpha + 2/(T_{\text{calib}} + 1)$. Otherwise, $\tilde{s}_{(\lfloor (T_{\text{calib}}+1)\alpha/2 \rfloor),i}$ and $\tilde{s}_{(\lceil (T_{\text{calib}}+1)(1-\alpha/2) \rceil),i}$ correspond respectively to some $\tilde{s}_{t,i}$ for some $t \in \mathcal{D}_{\text{calib}}$.

We apply arguments of exchangeability in the latter case. The new score $\tilde{s}_{T+1,i}$ is equally likely to fall into any of the $T_{\text{calib}} + 1$ intervals defined by the $(\tilde{s}_t)_{t \in \mathcal{D}_{\text{calib}}}$. More formally, by Assumption 3.1, and when scores are almost-surely distinct,

$$\mathbb{P}\left(\widetilde{s}_{T+1,i} < \widetilde{s}_{(1),i}\right) = \mathbb{P}\left(\widetilde{s}_{T+1,i} > \widetilde{s}_{(T_{\text{calib}}),i}\right) = \frac{1}{T_{\text{calib}} + 1}$$

and
$$\forall k \in [T_{\text{calib}} - 1],$$
 $\mathbb{P}(\widetilde{s}_{(k),i} < \widetilde{s}_{T+1,i} < \widetilde{s}_{(k+1),i}) = \frac{1}{T_{\text{calib}} + 1}$

605 Therefore, when scores are almost-surely distinct, the event of interest (11) rewrites

$$\left\{y_{T+1,i} \in \widetilde{C}_{i}(\boldsymbol{x}_{T+1})\right\} \stackrel{\text{a.s.}}{=} \left\{\widetilde{s}_{\left(\lfloor (T_{\text{calib}}+1)\alpha/2 \rfloor\right),i} < \widetilde{s}_{T+1,i} < \widetilde{s}_{\left(\lceil (T_{\text{calib}}+1)(1-\alpha/2) \rceil\right),i}\right\}$$

and has a probability

$$\begin{split} \mathbb{P}\big(y_{T+1,i} \in \widetilde{C}_i(\boldsymbol{x}_{T+1})\big) &= \frac{\left\lceil (T_{\text{calib}}+1)(1-\alpha/2) \right\rceil - \left\lfloor (T_{\text{calib}}+1)\alpha/2 \right\rfloor}{T_{\text{calib}}+1} \\ &\leqslant \frac{\left((T_{\text{calib}}+1)(1-\alpha/2)+1 \right) - \left((T_{\text{calib}}+1)\alpha/2 - 1 \right)}{T_{\text{calib}}+1} = 1 - \alpha + \frac{2}{T_{\text{calib}}+1} \,, \end{split}$$

as claimed.

We now prove that $\mathbb{P}(y_{T+1,i} \in \tilde{C}_i(\boldsymbol{x}_{T+1})) \ge 1 - \alpha$ whether or not scores are almost-surely distinct. To do so, we show below that

$$\forall k \in [T_{\text{calib}}], \qquad \mathbb{P}\big(\widetilde{s}_{T+1,i} \leqslant \widetilde{s}_{(k),i}\big) \geqslant \frac{k}{T_{\text{calib}}+1} \qquad \text{and} \qquad \mathbb{P}\big(\widetilde{s}_{T+1,i} < \widetilde{s}_{(k),i}\big) \leqslant \frac{k}{T_{\text{calib}}+1}, \tag{12}$$

so that, given the rewriting (11), we will end up with

$$\mathbb{P}\left(y_{T+1,i} \in \widetilde{C}_i(\boldsymbol{x}_{T+1})\right) \geqslant \frac{\left\lceil (T_{\text{calib}}+1)(1-\alpha/2) \right\rceil - \left\lfloor (T_{\text{calib}}+1)\alpha/2 \right\rfloor}{T_{\text{calib}}+1} \geqslant \frac{(T_{\text{calib}}+1)(1-\alpha/2) - (T_{\text{calib}}+1)\alpha/2}{T_{\text{calib}}+1} = 1-\alpha$$

It only remains to show (12). The event $\{\tilde{s}_{T+1,i} \leq \tilde{s}_{(k),i}\}$ is exactly the fact that $\tilde{s}_{T+1,i}$ is among the k smallest elements of the $(\tilde{s}_t)_{t \in \mathcal{D}_{calib} \cup \{T+1\}}$. By exchangeability, the probability of the latter event is at least $k/(T_{calib} + 1)$; it may be larger if several scores take the same value as the k-th smallest value. Similarly, the event $\{\tilde{s}_{T+1,i} < \tilde{s}_{(k),i}\}$ is exactly the fact that $\tilde{s}_{T+1,i}$ is among the k smallest elements of the $(\tilde{s}_t)_{t \in \mathcal{D}_{calib} \cup \{T+1\}}$ and that there are no ties at the k-th smallest value. Due to the additional no-tie condition, and by exchangeability, the probability of the latter event is at most $k/(T_{calib} + 1)$.

B. Proof of the efficiency result for fixed weights w (Theorem 3.7)

We first state some elementary properties of elliptical distributions.

Property B.1. The marginals of a spherical distribution are identically distributed. A spherical distribution with a first-order moment is centered: $\mathbb{E}[\boldsymbol{z}] = \boldsymbol{0}$. A spherical distribution with a second-order moment has a covariance matrix proportional to the identity: there exists $\sigma^2 \in [0, +\infty)$ such that $\mathbb{E}[\boldsymbol{z}\boldsymbol{z}^{\top}] = \sigma^2 \operatorname{Id}_k$.

Proof. The first property is proved by considering permutation matrices Γ . The second property holds because u = 0 is the only vector $u \in \mathbb{R}^k$ such that $\Gamma u = u$ for all orthogonal matrices (first consider permutation matrices to get that all components of u are equal). For the third property, denote by Σ the covariance matrix of z. Since it is symmetric (positive semi-definite), there exists an orthogonal matrix Γ and a vector $\lambda \in \mathbb{R}^k$ (with non-negative elements) such that $\Gamma \Sigma \Gamma = \text{diag}(\lambda)$. Now, $\Gamma^{\top} z$ has the same distribution as z, thus their covariance matrices are equal, which shows that $\Sigma = \text{diag}(\lambda)$. As marginals have the same distribution, we finally get $\Sigma = \sigma^2 \text{Id}_k$ for some $\sigma^2 \in [0, +\infty)$, which is actually positive except if the distribution of z is a Dirac at 0.

A slightly more advanced result provides the form of the characteristic function of an elliptical distribution. Its proof is based on first showing that characteristic functions of spherical distributions are exactly of the form $u \mapsto \phi(u^T u)$, which is consistent with the fact that spherical distributions are centered. Actually, it may be seen that ϕ is the characteristic function of the common distribution of the marginals of z.

Lemma B.2 (Kollo & von Rosen, 2005, Theorem 2.3.5). Consider a random variable following an elliptical distribution over \mathbb{R}^m , of the form $\mathbf{c} + M\mathbf{z}$, for a deterministic vector $\mathbf{c} \in \mathbb{R}^m$, a $m \times k$ matrix M such that MM^{\top} has rank k, and a random vector \mathbf{z} following a spherical distribution over \mathbb{R}^k . The characteristic function of $\mathbf{c} + M\mathbf{z}$ is of the form

$$\forall \boldsymbol{u} \in \mathbb{R}^{m}, \qquad \mathbb{E}\Big[\exp\big(\mathrm{i}\boldsymbol{u}^{\mathsf{T}}(\boldsymbol{c}+M\boldsymbol{z})\big)\Big] = \exp\big(\mathrm{i}\boldsymbol{u}^{\mathsf{T}}\boldsymbol{c}\big) \ \phi\big(\boldsymbol{u}^{\mathsf{T}}MM^{\mathsf{T}}\boldsymbol{u}\big),$$

for some function $\phi : \mathbb{R} \to \mathbb{C}$ that only depends on the distribution of z.

Lemma B.2 is instrumental in showing that the normalized marginals of (a linear transformation of) an elliptical distribution have comparable univariate distributions (that are homothetical), as stated next.

Lemma B.3. With the setting and the notation of Lemma B.2, let N be any $m \times m$ matrix and consider the random vector s = N(c + Mz). Let $\Lambda = NMM^{\top}N^{\top}$. There exists a random variable v, following a univariate distribution induced by the spherical distribution of z, such that

$$\forall i \in [m], \qquad s_i - \mathbb{E}[s_i] \stackrel{\text{(d)}}{=} \sqrt{\Lambda_{i,i}} v.$$

Proof. By Lemma B.2, the characteristic function of $s - \mathbb{E}[s]$ is $u \in \mathbb{R}^m \mapsto \phi(u^{\top} \Lambda u)$. Thus, the characteristic function of each $s_i - \mathbb{E}[s_i]$ is $u \in \mathbb{R} \mapsto \phi(\Lambda_{i,i}u^2)$. This shows the stated result, for a random variable v with characteristic function ϕ .

A final preliminary is result justifies that the matrix P_w introduced in the statement of Theorem 3.7 is well defined.

Lemma B.4. The matrices $H^{\top}H$ and $H^{\top}WH$ are $n \times n$ symmetric positive definite matrices, where W is itself a $m \times m$ symmetric positive definite matrix. Thus, these matrices are invertible.

Proof. The form (1) of H entails that $H^{\top}H = \text{Id}_n + H^{\top}_{\text{sub}}H_{\text{sub}}$, where $H^{\top}_{\text{sub}}H_{\text{sub}}$ is symmetric positive semi-definite. Thus, $H^{\top}H$ is symmetric positive definite. Given it is symmetric positive definite, the matrix W may be decomposed as $W = N^{\top}N$ for some $n \times n$ invertible matrix N. The matrix $H^{\top}WH = (NH)^{\top}NH$ is symmetric positive semi-definite. We show that it is even symmetric positive definite: $u^{\top}(NH)^{\top}NHu = 0$ is equivalent to the standard Euclidean norm of NHu being null, thus to $Hu = \mathbf{0}$ (as N is invertible); given the form (1) of H, we conclude that $u^{\top}(NH)^{\top}NHu = 0$ is equivalent to $u = \mathbf{0}$, which is the definition of $H^{\top}WH = (NH)^{\top}NH$ being definite.

We are now ready to prove Theorem 3.7, which we restate first.

Theorem 3.7. Let $w \in \mathbb{R}^m$ be a vector of positive numbers. Under Assumptions 3.1 and 3.6 (i.i.d. scores with elliptical distribution), the hierarchical version of SCP (Algorithm 2) run with $P = P_w$, where

$$P_{\boldsymbol{w}} \stackrel{\text{def}}{=} H \left(H^{\top} \operatorname{diag}(\boldsymbol{w}) H \right)^{-1} H^{\top} \operatorname{diag}(\boldsymbol{w}) \,, \tag{3}$$

provides prediction rectangles that are more efficient than the ones output by the plain multivariate version of SCP (Algorithm 1) in the following sense:

$$\mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell\big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big)^2\right] \leqslant \mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell\big(\widehat{C}_i(\boldsymbol{x}_{T+1})\big)^2\right].$$

Proof. The matrix P_w satisfies $P_w H = H$, thus, as in the beginning of the proof of Theorem 3.2, we have that for all $t \in \mathcal{D}_{\text{calib}}$,

$$\widetilde{\boldsymbol{s}}_t \stackrel{\text{def}}{=} \boldsymbol{y}_t - P_{\boldsymbol{w}} \widehat{\boldsymbol{y}}_t = P_{\boldsymbol{w}} (\boldsymbol{y}_t - \widehat{\boldsymbol{y}}_t) = P_{\boldsymbol{w}} \widehat{\boldsymbol{s}}_t.$$

0 We let, for all $t \in \mathcal{D}_{calib}$ and $i \in [m]$,

$$\widehat{\xi}_{t,i} = \widehat{s}_{t,i} - \mathbb{E}[\widehat{s}_{1,i}]$$
 and $\widetilde{\xi}_{t,i} = \widetilde{s}_{t,i} - \mathbb{E}[\widetilde{s}_{1,i}]$

By Assumption 3.1 (i.i.d. scores), for each $i \in [m]$, the univariate random variables $\hat{\xi}_{t,i}$, where $t \in \mathcal{D}_{\text{calib}}$ are i.i.d.; a similar statement holds for the $\tilde{\xi}_{t,i}$, where $t \in \mathcal{D}_{\text{calib}}$. By Assumption 3.6 and Lemma B.3, there exist a matrix Γ of the form $\Gamma = MM^{T}$ and a random variable v such that, for each $i \in [m]$,

$$\widehat{\xi}_{t,i} \stackrel{\text{(d)}}{=} \sqrt{\Gamma_{i,i}} v \quad \text{and} \quad \widetilde{\xi}_{t,i} \stackrel{\text{(d)}}{=} \sqrt{\Gamma'_{i,i}} v, \quad \text{where} \quad \Gamma' = P_{\boldsymbol{w}} \Gamma P_{\boldsymbol{w}}^{\top}.$$

The Let $(v_t)_{t \in \mathcal{D}_{calib}}$ be i.i.d. random variables with the same distribution as v. We conclude from the facts above that for each $i \in [m]$,

$$\begin{array}{ccc} & & & \\ 712 \\ & & \\ 713 \\ & & \\ 714 \end{array} & \left(\widehat{s}_{t,i}\right)_{t \in \mathcal{D}_{\text{calib}}} \stackrel{\text{(d)}}{=} \left(\mathbb{E}[\widehat{s}_{1,i}] + \sqrt{\Gamma_{i,i}} \ v_t\right)_{t \in \mathcal{D}_{\text{calib}}} & \text{and} & \left(\widetilde{s}_{t,i}\right)_{t \in \mathcal{D}_{\text{calib}}} \stackrel{\text{(d)}}{=} \left(\mathbb{E}[\widetilde{s}_{1,i}] + \sqrt{\Gamma_{i,i}'} \ v_t\right)_{t \in \mathcal{D}_{\text{calib}}}. \end{array}$$

The same equalities in distributions hold for the corresponding order statistics: for each $i \in [m]$,

$$\left(\widehat{s}_{(t),i}\right)_{1\leqslant t\leqslant T_{\text{calib}}} \stackrel{\text{(d)}}{=} \left(\mathbb{E}\left[\widehat{s}_{1,i}\right] + \sqrt{\Gamma_{i,i}} \ v_{(t)}\right)_{1\leqslant t\leqslant T_{\text{calib}}} \quad \text{and} \quad \left(\widetilde{s}_{(t),i}\right)_{1\leqslant t\leqslant T_{\text{calib}}} \stackrel{\text{(d)}}{=} \left(\mathbb{E}\left[\widetilde{s}_{1,i}\right] + \sqrt{\Gamma_{i,i}'} \ v_{(t)}\right)_{1\leqslant t\leqslant T_{\text{calib}}}.$$

By following the conventions of Section 2.2 and letting $v_{(0)} = -\infty$ and $v_{(T_{\text{calib}}+1)} = +\infty$, we even have these equalities in distribution over vectors indexed by $0 \le t \le T_{\text{calib}} + 1$: for each $i \in [m]$,

$$\left(\widehat{s}_{(t),i}\right)_{0\leqslant t\leqslant T_{\text{calib}}+1} \stackrel{\text{(d)}}{=} \left(\mathbb{E}\left[\widehat{s}_{1,i}\right] + \sqrt{\Gamma_{i,i}} \ v_{(t)}\right)_{0\leqslant t\leqslant T_{\text{calib}}+1} \text{ and } \left(\widetilde{s}_{(t),i}\right)_{0\leqslant t\leqslant T_{\text{calib}}+1} \stackrel{\text{(d)}}{=} \left(\mathbb{E}\left[\widetilde{s}_{1,i}\right] + \sqrt{\Gamma_{i,i}'} \ v_{(t)}\right)_{0\leqslant t\leqslant T_{\text{calib}}+1} \left(\widetilde{s}_{(t),i}\right)_{0\leqslant t\leqslant T_{\text{calib}}+1} \left(\widetilde{s}_{(t),i}\right)_{0\leqslant t\leqslant T_{\text{calib}}+1} \left(\widetilde{s}_{(t),i}\right)_{0\leqslant t\leqslant T_{\text{calib}}+1} \right)_{0\leqslant t\leqslant T_{\text{calib}}+1} \left(\widetilde{s}_{(t),i}\right)_{0\leqslant T_{\text{calib}}+1}$$

Now, for each $i \in [m]$, by design of Algorithms 1 and 2, the lengths of the intervals $\widehat{C}_i(\boldsymbol{x}_{T+1})$ and $\widetilde{C}_i(\boldsymbol{x}_{T+1})$ output equals

$$\begin{split} \ell\big(\widehat{C}_i(\boldsymbol{x}_{T+1})\big) &= \widehat{s}_{\left(\left\lceil (T_{\text{calib}}+1)(1-\alpha/2)\right\rceil\right),i} - \widehat{s}_{\left(\left\lfloor (T_{\text{calib}}+1)\alpha/2\right\rfloor\right),i} \\ \text{and} \qquad \ell\big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big) &= \widetilde{s}_{\left(\left\lceil (T_{\text{calib}}+1)(1-\alpha/2)\right\rceil\right),i} - \widetilde{s}_{\left(\left\lfloor (T_{\text{calib}}+1)\alpha/2\right\rfloor\right),i} \,. \end{split}$$

Thus, letting $L_{\alpha} = v_{\left(\lceil (T_{\text{calib}}+1)(1-\alpha/2) \rceil \right)} - v_{\left(\lfloor (T_{\text{calib}}+1)\alpha/2 \rfloor \right)}$, where $L_{\alpha} \ge 0$ a.s., we finally proved that for each $i \in [m]$,

$$\ell(\widehat{C}_i(\boldsymbol{x}_{T+1})) \stackrel{\text{(d)}}{=} \sqrt{\Gamma_{i,i}} L_{\alpha} \quad \text{and} \quad \ell(\widetilde{C}_i(\boldsymbol{x}_{T+1})) \stackrel{\text{(d)}}{=} \sqrt{\Gamma'_{i,i}} L_{\alpha}.$$

We showed so far that

$$\mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell(\widehat{C}_i(\boldsymbol{x}_{T+1}))^2\right] = \left(\sum_{i=1}^{m} w_i \Gamma_{i,i}\right) \mathbb{E}[L_{\alpha}^2] \quad \text{and} \quad \mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell(\widetilde{C}_i(\boldsymbol{x}_{T+1}))^2\right] = \left(\sum_{i=1}^{m} w_i \ \Gamma'_{i,i}\right) \mathbb{E}[L_{\alpha}^2],$$

where $\mathbb{E}[L_{\alpha}^2] \ge 0$ is possibly infinite (in which case the stated result holds). The proof is concluded in the case $\mathbb{E}[L_{\alpha}^2] < +\infty$ by noting that

$$\operatorname{Tr}(\operatorname{diag}(\boldsymbol{w}) \Gamma) = \sum_{i=1}^{m} w_i \Gamma_{i,i} \ge \sum_{i=1}^{m} w_i \Gamma'_{i,i} = \operatorname{Tr}(\operatorname{diag}(\boldsymbol{w}) \Gamma') = \operatorname{Tr}(\operatorname{diag}(\boldsymbol{w}) P_{\boldsymbol{w}} \Gamma P_{\boldsymbol{w}}^{\top}),$$

which is guaranteed by the lemma below with $W = \text{diag}(\boldsymbol{w})$, since $\Gamma = MM^{\top}$ for some $m \times k$ matrix.

The first part of Lemma B.5 is elementary. Its second part is inspired by Panagiotelis et al. (2021, Theorem 3.2), which is a result about using orthogonal projections in the $\|\cdot\|_W$ -norm to derive distance-reducing properties, and by trace-minimization results that are classic in the literature of forecast reconciliation (like Lemma 3.9 to be found in Appendix C). We however see this second part as a new result of our own. See Appendix D.1, and in particular, the comments after(15), for more details.

Lemma B.5. Fix a symmetric positive definite matrix W and consider the associated inner product and induced norm

$$oldsymbol{u},oldsymbol{u}'\in\mathbb{R}^m\longmapsto\langleoldsymbol{u},oldsymbol{u}'
angle_W=\sqrt{oldsymbol{u}^ op\!Woldsymbol{u}'}\qquad and\qquadoldsymbol{u}\in\mathbb{R}^m\longmapsto\|oldsymbol{u}\|_W\stackrel{
m def}{=}\sqrt{oldsymbol{u}^ op\!Woldsymbol{u}}$$
 .

Then, $P_W \stackrel{\text{def}}{=} H(H^\top W H)^{-1} H^\top W$ is the orthogonal projection onto Im(H) in the $\|\cdot\|_W$ -norm. Furthermore, for all $m \times k$ matrices M,

$$0 \leq \operatorname{Tr}(WP_W M M^{\top} P_W^{\top}) \leq \operatorname{Tr}(W M M^{\top}).$$

Proof. First, P_W is indeed a projection onto Im(H): namely, $P_W P_W = P_W$ and $P_W H = H$. To show that P_W is an orthogonal projection for the $\|\cdot\|_W$ -norm, it suffices to note that for all $u, u' \in \mathbb{R}^m$,

$$\langle P_W \boldsymbol{u}, \, \boldsymbol{u}' \rangle_W \stackrel{\text{def}}{=} (P_W \boldsymbol{u})^{\mathsf{T}} W \boldsymbol{u}' = \boldsymbol{u}^{\mathsf{T}} W P_W \boldsymbol{u}' \stackrel{\text{def}}{=} \langle \boldsymbol{u}, \, P_W \boldsymbol{u}' \rangle_W \,,$$

where we used that $P_W^{\top}W = WP_W$, given the closed-form expression of P_W .

Now, let z' be a standard Gaussian random k-vector: $z' \sim \mathcal{N}(\mathbf{0}, \mathrm{Id}_k)$. On the one hand, given the orthogonality proved for P_W and by a Pythagorean theorem,

$$|P_W M \boldsymbol{z}'||_W^2 \leqslant ||M \boldsymbol{z}'||_W^2 \quad \text{a.s.}$$
⁽¹³⁾

Now, by definition of the $\|\cdot\|_W$ -norm and by elementary properties of the trace,

$$\mathbb{E}\left[\|P_{W}M\boldsymbol{z}'\|_{W}^{2}\right] = \mathbb{E}\left[(P_{W}M\boldsymbol{z}')^{\mathsf{T}}WP_{W}M\boldsymbol{z}'\right] = \mathbb{E}\left[\operatorname{Tr}\left(WP_{W}M\boldsymbol{z}'(P_{W}M\boldsymbol{z}')^{\mathsf{T}}\right)\right]$$
$$= \operatorname{Tr}\left(WP_{W}M\underbrace{\mathbb{E}\left[\boldsymbol{z}'(\boldsymbol{z}')^{\mathsf{T}}\right]}_{=\operatorname{Id}_{k}}M^{\mathsf{T}}P_{W}^{\mathsf{T}}\right) = \operatorname{Tr}\left(WP_{W}MM^{\mathsf{T}}P_{W}^{\mathsf{T}}\right).$$

Similarly, $\mathbb{E}[||M\boldsymbol{z}'||_W^2] = \operatorname{Tr}(WMM^{\top}).$

The inequality (13) and the two equalities proved above conclude the proof.

C. Proof of the component-wise efficiency results (Theorem 3.10 and Corollary 3.11)

The proof of Theorem 3.10, which we restate below, is based on a key equality established in the proof of Theorem 3.7 and on a result that is central in the theory of forecast reconciliation, namely Lemma 3.9 (re-stated and re-proved at the end of this section).

We recall that we denoted by

$$\widetilde{C}_1^\star(oldsymbol{x}_{T+1}) imes \ldots imes \widetilde{C}_m^\star(oldsymbol{x}_{T+1})$$
 and $\widetilde{C}_1(oldsymbol{x}_{T+1}) imes \ldots imes \widetilde{C}_m(oldsymbol{x}_{T+1})$

the prediction rectangles output by the hierarchical version of SCP (Algorithm 2) run with $P_{\Sigma^{-1}} = H (H^{\top} \Sigma^{-1} H)^{-1} H^{\top} \Sigma^{-1}$ and any other choice of a projection matrix P onto Im(H), respectively.

Theorem 3.10. Under Assumptions 3.1, 3.6, and 3.8 (i.i.d. scores with elliptical distribution admitting a second-order moment), the hierarchical version of SCP (Algorithm 2) run with $P = P_{\Sigma^{-1}}$ provides prediction rectangles more efficient than with any other choice of a projection matrix onto Im(H):

$$\forall i \in [m], \quad \mathbb{E}\Big[\ell\big(\widetilde{C}_i^{\star}(\boldsymbol{x}_{T+1})\big)^2\Big] \leqslant \mathbb{E}\Big[\ell\big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big)^2\Big].$$

Proof. The proof of Theorem 3.7 did not rely on the existence of a second-order moment, i.e., of a covariance matrix Σ for the distribution of the scores \hat{s}_t . (It did not even rely on the existence of a first-order moment.)

When such a second-order moment exists, we may modify the proof of Theorem 3.7 in the following way, to obtain expected lengths depending on Σ . We also note that though we wrote the beginning of that proof for a specific projection matrix P_w onto Im(H), it holds for all projection matrices P onto Im(H), and even for all matrices P such that PH = H. Namely, when Algorithm 2 is run with any projection matrix P onto Im(H),

$$\mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell\big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big)^2\right] = \left(\sum_{i=1}^{m} w_i \, \Gamma'_{i,i}\right) \mathbb{E}\left[L_{\alpha}^2\right] = \operatorname{Tr}\left(\operatorname{diag}(\boldsymbol{w}) \ P \Gamma P^{\mathsf{T}}\right) \ \mathbb{E}\left[L_{\alpha}^2\right],$$

where $\Gamma = MM^{\top}$ for some matrix M such that scores \hat{s}_t have the same distribution as some c + Mz with z following some spherical distribution. In particular, Assumption 3.8 and Property B.1 impose that M is a $m \times m$ matrix and they entail that there exists $\sigma^2 > 0$ such that $\Sigma = \sigma^2 MM^{\top} = \sigma^2 \Gamma$.

Therefore, we actually have, when Algorithm 2 is run with any projection matrix P onto Im(H),

$$\mathbb{E}\left[\sum_{i=1}^{m} w_i \ \ell\big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big)^2\right] = \left(\sum_{i=1}^{m} w_i \ \Gamma'_{i,i}\right) \mathbb{E}\left[L_{\alpha}^2\right] = \operatorname{Tr}\left(\operatorname{diag}(\boldsymbol{w}) \ P\Sigma P^{\mathsf{T}}\right) \ \frac{\mathbb{E}\left[L_{\alpha}^2\right]}{\sigma^2} .$$

Lemma B.5 shows that $P_{\Sigma^{-1}}$ is a projection matrix onto Im(H). Lemma 3.9 below shows that for all projections P onto Im(H) and all positive vectors $w \in \mathbb{R}^m$,

$$\operatorname{Tr}(\operatorname{diag}(\boldsymbol{w}) P_{\Sigma^{-1}} \Sigma P_{\Sigma^{-1}}^{\top}) \leq \operatorname{Tr}(\operatorname{diag}(\boldsymbol{w}) P \Sigma P^{\top})$$

Collecting all elements, whether $\mathbb{E}[L^2_{\alpha}] = +\infty$ or $\mathbb{E}[L^2_{\alpha}] \in [0, +\infty)$, we proved so far that when Algorithm 2 is run with any projection matrix P onto $\mathrm{Im}(H)$ to output prediction intervals \tilde{C}_i ,

$$\forall \boldsymbol{w} \in (0, +\infty)^m, \qquad \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell \big(\widetilde{C}_i^{\star}(\boldsymbol{x}_{T+1})\big)^2\right] \leqslant \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell \big(\widetilde{C}_i(\boldsymbol{x}_{T+1})\big)^2\right]. \tag{14}$$

We obtain the claimed component-wise inequalities by taking $w_i = 1$ for one component i and letting $w_i \to 0$ for $j \neq i$. \Box

We now move on to the proof of Corollary 3.11.

 Corollary 3.11. In the setting and under the assumptions of Theorem 3.10, we also have

$$\forall i \in [m], \quad \mathbb{E}\left[\ell\left(\widetilde{C}_{i}^{\star}(\boldsymbol{x}_{T+1})\right)^{2}\right] \leq \mathbb{E}\left[\ell\left(\widehat{C}_{i}(\boldsymbol{x}_{T+1})\right)^{2}\right],$$

where the prediction rectangle $\widehat{C}_1(\boldsymbol{x}_{T+1}) \times \ldots \times \widehat{C}_m(\boldsymbol{x}_{T+1})$ is output by the plain multivariate version of SCP (Algorithm 1).

Proof. The result follows from Theorems 3.7 and 3.10 (which both hold under the stronger set of assumptions of Theo-842 rem 3.10). More precisely, for each $w \in (0, +\infty)^m$, denote by \widetilde{C}_i^w the prediction intervals output by Algorithm 2 run with 843 $P = P_w$. Theorem 3.7 ensures that

$$\forall \boldsymbol{w} \in (0, +\infty)^m, \qquad \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell \big(\widetilde{C}_i^{\boldsymbol{w}}(\boldsymbol{x}_{T+1}) \big)^2 \right] \leqslant \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell \big(\widehat{C}_i(\boldsymbol{x}_{T+1}) \big)^2 \right]$$

Equality (14) in the proof of Theorem 3.10 states that

$$\forall \boldsymbol{w} \in (0, +\infty)^m, \qquad \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell \big(\widetilde{C}_i^{\star}(\boldsymbol{x}_{T+1})\big)^2\right] \leqslant \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell \big(\widetilde{C}_i^{\boldsymbol{w}}(\boldsymbol{x}_{T+1})\big)^2\right].$$

Combining these two inequalities, we have

$$\forall \boldsymbol{w} \in (0, +\infty)^m, \qquad \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell \big(\widetilde{C}_i^{\star}(\boldsymbol{x}_{T+1}) \big)^2 \right] \leqslant \mathbb{E}\left[\sum_{i=1}^m w_i \ \ell \big(\widehat{C}_i(\boldsymbol{x}_{T+1}) \big)^2 \right],$$

and we conclude the proof with the same limit arguments as after (14) in the proof of Theorem 3.10.

The following lemma is a deep and central result in the theory of forecast reconciliation. First stated for $W = Id_m$ by Wickramasuriya et al. (2019), this result has since been extended to symmetric positive semi-definite matrices W in Panagiotelis et al. (2021) and Ando & Narita (2024). We provide a short and elementary proof, which may actually be seen as a simplification of the proof by Ando & Narita (2024), an article entirely devoted to proving Lemma 3.9. The latter article sees the minimization problem at hand as a constrained minimization problem (given how projections onto Im(H) may be written), thus introduced a Lagrangian and discussed Karush-Kuhn-Tucker conditions to solve it.

Lemma 3.9 (Minimum-Trace projection). Let W and Σ be two symmetric $m \times m$ matrices, where W positive semi-definite and Σ is positive definite. Then, for all projection matrices P onto Im(H),

$$\operatorname{Tr}(WP_{\Sigma^{-1}}\Sigma P_{\Sigma^{-1}}^{\top}) \leq \operatorname{Tr}(WP\Sigma P^{\top}).$$

Proof. We first show that projection matrices P onto Im(H) are exactly the matrices of the form HG, where G is a $n \times m$ matrix such that $GH = Id_n$. Indeed, such a matrix HG satisfies HGHG = HG and HGH = H, which characterizes projections onto Im(H). Conversely, fix a projection P onto Im(H) and a basis u_1, \ldots, u_m of \mathbb{R}^m : each Pu_i belongs to Im(H), thus is of the form Hg_i for some $g_i \in \mathbb{R}^n$. Denote by G the $n \times m$ matrix with columns given by g_1, \ldots, g_m . By linearity of P and given that u_1, \ldots, u_m is a basis, we have P = HG. We denote by h_1, \ldots, h_n the columns of the $m \times n$ structural matrix H. Since P is a projection onto Im(H), we have in particular $Ph_i = h_i$ for all $i \in [n]$, or put differently, PH = H. Substituting P = HG and multiplying both sides by H^{T} , we proved so far that $H^{\mathsf{T}}HGH = H^{\mathsf{T}}H$, where (see Lemma B.4), the matrix $H^{\top}H$ is invertible. All in all, we thus proved $GH = \mathrm{Id}_n$.

880 Given the characterization above, the projection matrices P onto Im(H) are also exactly the matrices of the form

$$P = P_{\Sigma^{-1}} + HA = H\left(\left(H^{\mathsf{T}}\Sigma^{-1}H\right)^{-1}H^{\mathsf{T}}\Sigma^{-1} + A\right), \qquad \text{for } n \times m \text{ matrices } A \text{ such that} \qquad AH = [0]_n,$$

where $[0]_n$ denotes the $n \times n$ null matrix. Keeping in mind that Σ and Σ^{-1} are symmetric, this decomposition entails that

$$\begin{split} WP\Sigma P^{\mathsf{T}} &= W\left(H\left(H^{\mathsf{T}}\Sigma^{-1}H\right)^{-1}H^{\mathsf{T}}\Sigma^{-1}\right)\Sigma\left(\Sigma^{-1}H\left(H^{\mathsf{T}}H\right)^{-1}H^{\mathsf{T}}\right) \\ &+ W\left(HA\right)\Sigma\left(\Sigma^{-1}H\left(H^{\mathsf{T}}\Sigma^{-1}H\right)^{-1}H^{\mathsf{T}}\right) \\ &+ W\left(H\left(H^{\mathsf{T}}\Sigma^{-1}H\right)^{-1}H^{\mathsf{T}}\Sigma^{-1}\right)\Sigma\left(A^{\mathsf{T}}H^{\mathsf{T}}\right) \\ &+ W\left(HA\right)\Sigma\left(HA\right)^{\mathsf{T}}. \end{split}$$

The second term in the decomposition simplifies into

$$W(HA)\Sigma\left(\Sigma^{-1}H\left(H^{\mathsf{T}}\Sigma^{-1}H\right)^{-1}H^{\mathsf{T}}\right) = WH\overbrace{AH}^{-\overset{(0)n}{(AH)}}\left(H^{\mathsf{T}}\Sigma^{-1}H\right)^{-1}H^{\mathsf{T}} = [0]_{m}.$$

-[0]

Similarly, the third term is also null, due to the term $H^{\top}\Sigma^{-1}\Sigma A^{\top} = (AH)^{\top}$. The proof is concluded by noting that for all matrices A, the trace of the fourth term in the decomposition of $WP\Sigma P^{\top}$ is non-negative. Indeed, given that W and Σ are positive semi-definite, we may write them as $W = MM^{\top}$ and $\Sigma = NN^{\top}$ for $m \times m$ matrices M and N. Then, together with elementary properties of the trace,

$$\begin{aligned} \operatorname{Tr} \Big(W \left(HA \right) \Sigma \left(HA \right)^{\top} \Big) &= \operatorname{Tr} \Big(M M^{\top} (HA) N N^{\top} (HA)^{\top} \Big) = \operatorname{Tr} \Big(M^{\top} (HA) N N^{\top} (HA)^{\top} M \Big) \\ &= \operatorname{Tr} \Big(\left(M^{\top} (HA) N \right) \left(M^{\top} (HA) N \right)^{\top} \Big) \geqslant 0 \,, \end{aligned}$$

given that the trace of a symmetric positive semi-definite matrix is non-negative.

D. Forecast reconciliation: review and connections made

For a complete review on the forecast reconciliation literature, we refer the reader to Athanasopoulos et al. (2024) and only
 provide a brief overview below.

At the origins, forecasts in the hierarchical setting were conducted using a single-level approach (most notably, in the bottom-up or top-down fashion), i.e., by choosing a level of the hierarchy (typically, either the bottom level or the top level) to generate forecasts, and then, by propagating these forecasts (typically in a linear fashion). A notable pitfall of the single-level approaches is that potentially valuable information from all other levels are ignored. To overcome this issue, the concept of forecast reconciliation was introduced by Athanasopoulos et al. (2009) and Hyndman et al. (2011): the idea is to combine forecasts from different levels of aggregation through linear combinations. Recently, developments were made in reconciliation through projections (Wickramasuriya et al., 2019, Panagiotelis et al., 2021), which we review and detail in the next section.

Probabilistic hierarchical forecasting and reconciliation is an emerging field. Notable works include the one by Wickrama suriya (2024), which studied probabilistic forecast reconciliation for Gaussian distributions, while Panagiotelis et al. (2023)
 provided reconciled forecasts based on the minimization of a probabilistic score through gradient descent. However, we did
 not leverage results from this literature for our own probabilistic approach.

926 927 **D.1. Review of forecast reconciliation through projections**

We summarize and review the approach followed by Hyndman et al. (2011), Wickramasuriya et al. (2019), and Panagiotelis et al. (2021).

⁹³⁰The setting is the one described in Section 2, with stochastic observations following some hierarchical structure $y = Hy_{1:n}$; ⁹³²features are possibly available. Initial point forecasts \hat{y} are provided by some regression method A; these forecasts are ⁹³³possibly incoherent, i.e., do not belong to Im(H). The goal of forecast reconciliation is to leverage the hierarchical structure ⁹³⁴to improve the point forecasts.

A typical assumption made in this literature is that the point forecasts \hat{y} are unbiased, or, put differently, that the forecasting errors $\hat{s} = y - \hat{y}$ are centered. A natural performance criterion then is the mean-square error [MSE]: letting $\|\cdot\|_2$ denote the Euclidean norm and Σ the covariance matrix of $\hat{s} = y - \hat{y}$,

$$\mathsf{MSE}(\widehat{\boldsymbol{y}}, \boldsymbol{y}) \stackrel{\mathrm{def}}{=} \mathbb{E}[\|\widehat{\boldsymbol{s}}\|_{2}^{2}] = \mathbb{E}[\widehat{\boldsymbol{s}}^{\top}\widehat{\boldsymbol{s}}] = \mathbb{E}[\mathrm{Tr}(\widehat{\boldsymbol{s}}\widehat{\boldsymbol{s}}^{\top})] = \mathrm{Tr}(\mathbb{E}[\widehat{\boldsymbol{s}}\widehat{\boldsymbol{s}}^{\top}]) \stackrel{\mathrm{def}}{=} \mathrm{Tr}(\Sigma).$$

The equalities above may be generalized to W-norms (as defined in Lemma B.5), where W is a symmetric definite positive matrix:

$$\mathsf{MSE}(\widehat{\boldsymbol{y}}, \boldsymbol{y}, W) \stackrel{\text{def}}{=} \mathbb{E}[\|\widehat{\boldsymbol{s}}\|_W^2] = \mathbb{E}[\widehat{\boldsymbol{s}}^\top W \widehat{\boldsymbol{s}}] = \mathbb{E}[\mathrm{Tr}(W \widehat{\boldsymbol{s}} \widehat{\boldsymbol{s}}^\top)] = \mathrm{Tr}(W \Sigma).$$

Natural improvements of the unbiased point forecasts are exactly given by projections thereof onto Im(H), as justified below in Lemma D.1. Let P be a projection onto Im(H) and denote $\tilde{y} = P\hat{y}$. By linearity of a projection, the point forecasts \tilde{y} are also unbiased. Since observations are coherent, we have

$$\boldsymbol{y} - \widetilde{\boldsymbol{y}} \stackrel{\text{def}}{=} \boldsymbol{y} - P\widehat{\boldsymbol{y}} = P(\boldsymbol{y} - \widehat{\boldsymbol{y}}) = P\widehat{\boldsymbol{s}} \stackrel{\text{def}}{=} \widetilde{\boldsymbol{s}}.$$

The mean-squared error of \tilde{y} in W-norm thus equals

$$\mathsf{MSE}(\widetilde{\boldsymbol{y}}, \boldsymbol{y}, W) = \mathbb{E}\Big[\mathrm{Tr}\big(W\widetilde{\boldsymbol{s}}\,\widetilde{\boldsymbol{s}}^{\mathsf{T}}\big)\Big] = \mathrm{Tr}\Big(W\,P\,\mathbb{E}\big[\widehat{\boldsymbol{s}}\,\widehat{\boldsymbol{s}}^{\mathsf{T}}\big]P^{\mathsf{T}}\Big) = \mathrm{Tr}\big(W\,P\Sigma P^{\mathsf{T}}\big)\,.$$

Actually, the formula above holds more generally for all matrices P such that PH = H.

Optimal unbiased point forecasts in the sense of the mean-square error thus exactly correspond to minimizing $\text{Tr}(W P \Sigma P^{T})$, a problem that we discuss below. Before we do so, we justify why (only) projections onto Im(H) are considered.

Why (only) projections onto Im(H) are considered. This follows from the lemma below, given that the literature of forecast reconciliation considers, implicitly or explicitly, two restrictions: that forecasts should be unbiased; that improved forecasts should be obtained by linear combinations of the original forecasts (and be coherent, of course).

Lemma D.1 (Hyndman et al., 2011). Assume that the point forecasts \hat{y} are unbiased. Let M be a $m \times m$ matrix taking values in the coherent subspace Im(H). Then the linear combinations $\tilde{y} = M\hat{y}$ are unbiased if and only if M is a projection onto Im(H).

Proof. Being unbiased means the following in Hyndman et al. (2011): we denote by $m = H\beta$ the expectation of y, i.e., $\mathbb{E}[y] = m = H\beta$, and assume that the model is rich enough so that all values of $\beta \in \mathbb{R}^n$, i.e., all values of $m \in \text{Im}(H)$, may be obtained when the specifications of the model vary.

That $\widetilde{y} = M \widehat{y}$ is unbiased thus corresponds to the equalities

$$\forall \boldsymbol{\beta} \in \mathbb{R}^n, \quad MH\boldsymbol{\beta} = H\boldsymbol{\beta}, \quad \text{i.e.,} \quad MH = H.$$

Now, the proof of Lemma 3.9 in Appendix C shows that since M takes values in Im(H), it is of the form M = HG for some $n \times m$ matrix G. The equality MH = H may be rewritten as HGH = H. Again as in the proof of Lemma 3.9, by multiplying both sides of this equality by $(H^{\top}H)^{-1}H^{\top}$, we obtain $GH = \text{Id}_n$, which yields $M^2 = HGHG = HG = M$. Thus, M is indeed a projection onto Im(H).

Trace optimization, part 1: known covariance matrix. As explained above, original unbiased forecasts \hat{y} and their (still unbiased, linear) transformations $\tilde{y} = P\hat{y}$, where P is a projection onto Im(H) may be compared through their mean-squared errors in W-norm:

$$MSE(\widehat{\boldsymbol{y}}, \boldsymbol{y}, W) = Tr(W\Sigma)$$
 vs. $MSE(\widetilde{\boldsymbol{y}}, \boldsymbol{y}, W) = Tr(WP\Sigma P^{\top})$.

This consideration leads to the central result in forecast reconciliation: the optimality of the so-called Minimum Trace reconciliation method from Wickramasuriya et al. (2019), formally re-stated below as Lemma 3.9. This method consists of projecting according to $P_{\Sigma^{-1}} \stackrel{\text{def}}{=} H (H^{\top} \Sigma^{-1} H)^{-1} H^{\top} \Sigma^{-1}$, where we recall that Σ is the (unknown) covariance matrix of the forecast errors \hat{s} . Of course, this theoretically optimal method must be turned into a practical method, e.g., by replacing Σ in the formula above by some empirical estimate. Lemma 3.9 was originally stated by Wickramasuriya et al. (2019) in the case $W = \text{Id}_m$, and later extended to symmetric positive semi-definite matrices W by Panagiotelis et al. (2021) and Ando & Narita (2024). As discussed in Appendix C, we provide a more elementary proof.

⁹⁹³ **Lemma 3.9** (Minimum-Trace projection). Let W and Σ be two symmetric $m \times m$ matrices, where W positive semi-definite and Σ is positive definite. Then, for all projection matrices P onto Im(H),

$$\operatorname{Tr}(WP_{\Sigma^{-1}}\Sigma P_{\Sigma^{-1}}^{\top}) \leq \operatorname{Tr}(WP\Sigma P^{\top}).$$

Trace optimization, part 2: a more practical approach. The drawback with the approach above is that it relies on the knowledge of the covariance matrix Σ , but its advantage is that it holds for all weight matrices W. We now show how to exchange the roles of W and Σ , and get a trace-reduction result for a given weight matrix W but for all possible covariance matrices Γ , i.e., symmetric positive semi-definite matrices.

This result is inspired from Panagiotelis et al. (2021), who recommend to use the orthogonal projection in *W*-norm, whose closed-form expression (see Lemma B.5) reads $P_W \stackrel{\text{def}}{=} H(H^\top W H)^{-1} H^\top W$. A Pythagorean theorem ensures that, for all point forecasts \hat{y} and (coherent) observations y,

$$\left\| \boldsymbol{y} - P_W \widehat{\boldsymbol{y}} \right\|_W^2 = \left\| P_W (\boldsymbol{y} - \widehat{\boldsymbol{y}}) \right\|_W^2 \leqslant \left\| (\boldsymbol{y} - \widehat{\boldsymbol{y}}) \right\|_W^2$$
 a.s.,

thus, by taking expectations,

996 997

1002

1006 1007

1009

1014

1019

1043 1044

$$\operatorname{Tr}(W P_W \Sigma P_W^{\top}) = \operatorname{MSE}(P_W \widehat{y}, y, W) \leq \operatorname{MSE}(\widehat{y}, y, W) = \operatorname{Tr}(W \Sigma)$$

The equality above holds no matter the specific value of the covariance matrix Σ , which corresponds to the following trace-reduction inequality, stated as the second part of Lemma B.5: for all symmetric positive semi-definite matrices Γ ,

$$0 \leqslant \operatorname{Tr}(WP_W \Gamma P_W^{\dagger}) \leqslant \operatorname{Tr}(W\Gamma) . \tag{15}$$

The inequality above (i.e., Lemma B.5) is a result of our own though it was inspired by both Lemma 3.9 and the approach by Panagiotelis et al. (2021) relying on P_W -projections.

D.2. How we leveraged and transferred these results (and why it was not immediate)

On the unnecessity of unbiasedness. As we made clear several times in Section D.1, a key assumption in the literature of forecast reconciliation is that point forecasts are unbiased, or put differently, that the forecasting errors \hat{s} are centered.

This is in sharp contrast with the non-conformity scores \hat{s} considered in this article, which we do not want (nor need) to assume are centered. None of Assumptions 3.1–3.6–3.8 are about this. We rather assume that these scores follows a so-called elliptic distribution, with possibly a non-null expectation. Elliptic distributions were considered, not in the literature of reconciliation of point forecasts but of probabilistic forecasts, see Panagiotelis et al. (2023). Now, the proof of Theorem 3.7 in Appendix B reveals that our construction of prediction rectangles is such that the length of the *i*-th defining interval is given by

$$\ell(C_i(\boldsymbol{x}_{T+1})) = \widehat{s}_{\left(\lceil (T_{\text{calib}}+1)(1-\alpha/2)\rceil\right), i} - \widehat{s}_{\left(\lfloor (T_{\text{calib}}+1)\alpha/2\rfloor\right), i}$$

Non-null expectations of the underlying elliptic distribution cancel out in the above equation, hence the unnecessity of an assumption of unbiasedness. The cancellation is only possible because we considered signed non-conformity scores (which is unusual in the literature of conformal prediction).

On the component-wise objective ($\star\star$). As we detail in Appendix E (see the comments before Theorem E.3), the theory provided in this article is only worth being detailed because we do not target joint-coverage guarantees but component-wise coverage guarantees. We had to find out a component-wise efficiency objective that we could handle. With the literature of forecast reconciliation in mind, we somehow had to build an intuition a such an efficiency criterion.

1040 The proof of Theorem 3.7 in Appendix B explains how we could relate our (component-wise) small-length objective ($\star\star$), 1041 namely,

minimizing
$$\mathbb{E}\left[\sum_{i=1}^{m} w_i \ell (C_i(\boldsymbol{x}_{T+1}))^2\right],$$
 (16)

1045 to problems of the form

1046

minimizing
$$\operatorname{Tr}(\operatorname{diag}(\boldsymbol{w}) P \Gamma P),$$
 (17)

for some symmetric positive semi-definite matrix Γ , so as to leverage inequality (15), which is of our own. The proof of Theorem 3.10 reveals that when non-conformity scores have a bounded second-order moment, the matrix Γ is proportional to their covariance matrix Σ , which opened the avenue of the Minimum Trace approaches of Lemma 3.9.

Summary of the challenges overcome. In a nutshell, the main challenge overcome was to relate the two minimization problems (16) and (17), and in the first place, state suitably the efficiency criterion (16). The main tools used were to resort to signed non-conformity scores, which are not necessarily unbiased, and to exploit properties of elliptic distributions, in terms of stability of the shapes of these distributions under certain affine transformations.

7 E. Multi-target split conformal prediction: Reminder and extension

Conformal prediction is a framework initially introduced to quantity the uncertainty around univariate targets thanks to univariate non-conformity scores. In this section, we focus on extensions of conformal prediction to multivariate targets as introduced by Johnstone & Cox (2021) and Messoudi et al. (2022). The key idea of these methods is to consider A-norms of non-conformity scores, where A a data-based definite positive matrix designed to capture the potential multivariate dependencies of the targets. (The choice of A is the key for efficient prediction regions in practice; see the statement of Lemma B.5 for a reminder of the definition of A-norms.) In a nutshell, and as formally detailed below, the consideration of data-based norms effectively matches vector predictions to scalar non-conformity scores.

1066 1067 **E.1. Overview of existing results**

We first state the objectives and then review the methodology followed. We consider multivariate data (features $x_t \in \mathbb{R}^d$ and observations $y_t \in \mathbb{R}^m$) but with no specific hierarchical structure—unlike in Section E.2 below.

We now denote the prediction regions by $E(\mathbf{x})$ because they will typically be given by ellipsoids. (Messoudi et al., 2022 empirically illustrated that ellipsoidal predictive regions are more efficient than hyper-rectangular ones in terms of volumes.)

1074 **Objectives.** The cited references are interested in joint coverage guarantees and replace the component-wise coverage 1075 objective (*) by the design of prediction regions $E : \mathbf{x} \mapsto E(\mathbf{x}) \subseteq \mathbb{R}^m$ such that

$$\mathbb{P}(\boldsymbol{y}_{T+1} \in E(\boldsymbol{x}_{T+1})) \approx 1 - \alpha, \qquad (\diamondsuit)$$

where the probability \mathbb{P} is with respect to both $(\boldsymbol{x}_{T+1}, \boldsymbol{y}_{T+1})$ and $(\boldsymbol{x}_t, \boldsymbol{y}_t)_{1 \leq t \leq T}$.

The secondary objective of ensuring that the prediction sets output are efficient, i.e., as small as possible. We could state it through volumes, as given by the Lebesgue measure \mathfrak{L} over \mathbb{R}^m , and replace objective (**) by

minimizing
$$\mathbb{E}\left[\mathfrak{L}(E(\boldsymbol{x}_{T+1}))\right],$$
 (\diamondsuit)

where again, the expectation is with respect to both $(\boldsymbol{x}_{T+1}, \boldsymbol{y}_{T+1})$ and $(\boldsymbol{x}_t, \boldsymbol{y}_t)_{1 \leq t \leq T}$. In Theorem E.3, we are actually able to prove an even stronger result of uniform domination: a prediction region E is uniformly more efficient than a prediction region E' if $E(\boldsymbol{x}) \subseteq E'(\boldsymbol{x})$ for all $\boldsymbol{x} \in \mathbb{R}^d$.

1089 **Methodology and algorithm.** We actually present a generalization of the methodology and algorithm considered 1090 by Johnstone & Cox (2021), in the spirit of Algorithm 3. The proof of the joint coverage result below, Theorem E.2, only 1091 relies on Assumption E.1, which holds as long as the matrix A is chosen based on data of \mathcal{D}_{train} and \mathcal{D}_{est} only (just as was 1092 the case in Section 3.4).

The exact procedure is stated in Algorithm 4. A regressor $\hat{\mu}$ is built based on the data of the train set $\mathcal{D}_{\text{train}}$ and on a regression algorithm \mathcal{A} . Then, based on data $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{estim}}$, but preferably only on $\hat{\mu}$ and on data from $\mathcal{D}_{\text{estim}}$, some symmetric definite positive matrix A is computed in some black-box fashion; this is why Algorithm 4 takes a procedure \mathcal{E} as input. This matrix A is used to define a weighted norm (see Lemma B.5):

1098

1099

1077

$$oldsymbol{u} \in \mathbb{R}^m \longmapsto \|oldsymbol{u}\|_A \stackrel{ ext{def}}{=} \sqrt{oldsymbol{u}^ op Aoldsymbol{u}}$$

Conformal Prediction for Hierarchical Data



7: return $\check{E}(\boldsymbol{x}_{T+1})$

The multivariate prediction errors $y_t - \hat{\mu}(x_t)$ for $t \in \mathcal{D}_{calib}$ are then transformed into the univariate non-conformity scores through *A*–norms:

$$\widehat{s}_t = \| \boldsymbol{y}_t - \widehat{\boldsymbol{y}}_t \|_A$$
.

Some threshold $\check{q}_{1-\alpha}$ is determined based on the empirical quantiles of the series of these scores, and finally, the following prediction regions are output, which are ellipsoids:

$$\check{E}(\,\cdot\,) = \left\{ \boldsymbol{y} \in \mathbb{R}^m : \left\| \boldsymbol{y} - \widehat{\boldsymbol{\mu}}(\,\cdot\,) \right\|_A \leqslant \check{q}_{1-\alpha} \right\}.$$

For instance, Johnstone & Cox (2021) use the so-called Mahalanobis distance, which corresponds to taking A as the inverse of the covariance matrix of the forecasting errors (in the same spirit as in Algorithm 3).

Coverage guarantees. They rely on an i.i.d. assumption on the non-conformity scores (just as for Theorem 3.2 and as, more generally, in the literature of conformal prediction).

Assumption E.1. The non-conformity scores $\check{s}_t = \| \boldsymbol{y}_t - \hat{\boldsymbol{\mu}}(\boldsymbol{x}_t) \|_A$ are i.i.d. for $t \in \mathcal{D}_{\text{calib}} \cup \{T+1\}$. This is in particular the case when data $(\boldsymbol{x}_t, \boldsymbol{y}_t)_{1 \le t \le T+1}$ is i.i.d.

Theorem E.2. Fix $\alpha \in (0, 1)$. Algorithm 4, used with any regression algorithm \mathcal{A} and any estimation procedure \mathcal{E} , ensures that whenever Assumption E.1 (i.i.d. scores) holds,

$$\mathbb{P}(\boldsymbol{y}_{T+1} \in \check{E}(\boldsymbol{x}_{T+1})) \ge 1 - \alpha$$

In addition, if the non-conformity scores $(\check{s}_t)_{t \in \mathcal{D}_{calib} \cup \{T+1\}}$ are almost surely distinct, then

$$\forall i \in [m], \mathbb{P} \big(\boldsymbol{y}_{T+1} \in \check{E}(\boldsymbol{x}_{T+1}) \big) \leqslant 1 - \alpha + \frac{1}{T_{\text{calib}} + 1}.$$

Proof. We first note that by definition,

$$\left\{\boldsymbol{y}_{T+1} \in \check{E}(\boldsymbol{x}_{T+1})\right\} = \left\{\check{s}_{T+1} \leqslant \check{s}_{\left(\left\lceil (T_{\text{calib}}+1)(1-\alpha)\right\rceil\right)}\right\}.$$

The rest of the proof consists of the same classical arguments that were already detailed in the proof of Theorem 3.2 in Appendix A.

E.2. Extension to hierarchical data

We now assume that the observations y_t follow a hierarchical structure, as in (1). We adapt Algorithm 4 into Algorithm 5 in the same way we obtained Algorithm 2 from Algorithm 1, by merely adding a projection step right before computing the non-conformity scores on $\mathcal{D}_{\text{calib}}$. We resort to the orthogonal projection matrix onto Im(H) in $\|\cdot\|_A$ -norm, whose closed-form expression reads (see Lemma B.5)

$$P_A = H(H^{\mathsf{T}}AH)^{-1}H^{\mathsf{T}}A$$

1155	Algorithm 5 Hierarchical SCP through data-based norms
1156	Parameters: confidence level $1 - \alpha$; regression algorithm \mathcal{A} ; partition of $[T]$ into subsets $\mathcal{D}_{\text{train}}$, $\mathcal{D}_{\text{estim}}$ and $\mathcal{D}_{\text{calib}}$ of respective
1157	cardinalities T_{train} , T_{estim} and T_{calib} ; estimation procedure \mathcal{E} of the matrix used to define the norm
1158	1: Build the regressor $\widehat{\mu}(\cdot) = \mathcal{A}((x_t, y_t)_{t \in \mathcal{D}_{\text{train}}})$
1159	2: Compute a symmetric definite positive matrix $A = \mathcal{E}((\boldsymbol{x}_t, \boldsymbol{y}_t)_{t \in \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{estim}}})$
1161	3: Let $P_A = H(H^{\top}AH)^{-1}H^{\top}A$ and consider $\mathring{\mu}(\cdot) = P_A\widehat{\mu}(\cdot)$
1162	4: for $t \in \mathcal{D}_{\text{calib}}$ let $\mathring{m{y}}_t = \mathring{m{\mu}}(m{x}_t)$ and $\mathring{s}_t = \ m{y}_t - \mathring{m{y}}_t\ _A$
1162	5: Order the $(\mathring{s}_t)_{t \in \mathcal{D}_{\text{calib}}}$ into $\mathring{s}_{(1)} \leq \ldots \leq \mathring{s}_{(T_{\text{calib}})}$ and define $\mathring{s}_{(0)} = 0$ and $\mathring{s}_{(T_{\text{calib}}+1)} = +\infty$
1164	6: Let $\mathring{q}_{1-\alpha} = \mathring{s}_{\left(\lceil (T_{\text{calib}}+1)(1-\alpha) \rceil \right)}$
1165	7: Set $\mathring{E}(\cdot) = \left\{ \boldsymbol{y} \in \mathbb{R}^m : \left\ \boldsymbol{y} - \widehat{\boldsymbol{\mu}}(\cdot) \right\ _A \leqslant \mathring{q}_{1-\alpha} \right\}$
1166	8: return $\mathring{E}(\boldsymbol{x}_{T+1})$

1169 1170 Instead of considering the regressor $\hat{\mu}$, which may yield point estimates not abiding by the hierarchical constraints, we resort 1171 to $\mathring{\mu} = P_A \hat{\mu}$. The rest of the procedure is similar to the previous section.

¹¹⁷² We summarize the adaptation in Algorithm 5. The statement only differs from Algorithm 4 by the addition of the blue line. ¹¹⁷³ We denote by \mathring{s} and \mathring{E} the scores and prediction ellipsoids obtained after the projection step, to distinguish them from the ¹¹⁷⁴ ones obtained by Algorithm 4 without the projection step.

We may prove the following result, which shows that the objectives (\diamond) and ($\diamond \diamond$) are met. Its proof is straightforward. Put differently, there would have been no challenge in providing a theory of efficient conformal prediction for hierarchical data under a joint-coverage objective (\diamond). This was not the case at all for component-wise coverage objectives, as the tools of forecast reconciliation (like the projections step by P_A) are not component-wise tools. The proof of Theorem E.3 actually emphasizes the complexity of results such as Theorems 3.10–3.7 and Corollary 3.11.

1181 **Theorem E.3.** Fix $\alpha \in (0,1)$. Algorithm 5, used with any regression algorithm \mathcal{A} and any estimation procedure \mathcal{E} , 1182 guarantees the same coverage results as in Theorem E.2 whenever Assumption E.1 (i.i.d. scores) holds.

¹¹⁸³ ¹¹⁸⁴ ¹¹⁸⁴ ¹¹⁸⁵ In addition, the prediction ellipsoids \mathring{E} output by Algorithm 5 are uniformly more efficient than the prediction ellipsoids \check{E} ¹¹⁸⁵ output by Algorithm 4:

$$\check{E}(\boldsymbol{x}_{T+1}) \subseteq \check{E}(\boldsymbol{x}_{T+1})$$
 a.s.

Proof. The proof of the coverage guarantees is similar to the one in Theorem E.2. For the uniform efficiency part, we first note by a Pythagorean theorem, and since observations are coherent, $\dot{s}_t \leq \check{s}_t$ for all $t \in \mathcal{D}_{calib}$. Thus, in particular

$$\mathring{s}_{\left(\left\lceil (T_{\text{calib}}+1)(1-\alpha)\right\rceil\right)} \leqslant \mathring{s}_{\left(\left\lceil (T_{\text{calib}}+1)(1-\alpha)\right\rceil\right)},$$

hence, the stated inclusion, by construction of the predictive regions.

1194 1195 **F. Full details for the simulations: settings, methodology, results**

In this appendix, we provide the full details on the specifications and of the results of the numerical experiments summarized
 in the main body of the article.

11991200F.1. Artificially generated data

1186 1187

1193

1201 The objective of the experiments on synthetic data is to replicate the behavior of real data while controlling the number of 1202 observations available and the difficulty of the forecasting task.

1204 F.1.1. DATA GENERATION

Data generation: initial draw of the parameters. The structural matrix H considered in this example is the right one in Figure 1; we copy it in Figure 3 for the convenience of the reader. Therefore, there are m = 8 nodes in the hierarchy, with n = 5 nodes at the most disaggregated levels. For each given run, we first pick at random a function $f : \mathbb{R}^3 \to \mathbb{R}^5$ and a covariance matrix $A^T A$. We do so as explained later in this description.



Figure 3. The structural matrix H considered in the numerical experiments with artificially generated data.

Data generation: draw of *T*-sample. Then, given *H*, *f*, and *A*, we draw a *T*-sample $(x_t, y_t)_{1 \le t \le T}$ of data as follows. First, the features $x_t \in \mathbb{R}^3$ are drawn i.i.d. according to a Gaussian distribution:

$$\boldsymbol{x}_t = \begin{bmatrix} x_{t,1} \\ x_{t,2} \\ x_{t,3} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}, \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right).$$

Next, the observations $y_{t,1:5} \in \mathbb{R}^5$ at the most disaggregated level are generated i.i.d. according to the following additive model:

$$\boldsymbol{y}_{t,1:5} = f(\boldsymbol{x}_t) + \boldsymbol{\varepsilon}_t, \quad \text{where} \quad \boldsymbol{\varepsilon}_t \sim \mathcal{N}\left(\begin{bmatrix} 10\\ \vdots\\ 10 \end{bmatrix}, A^{\mathsf{T}}A \right).$$
 (18)

1230 The complete vectors of observations are finally given by $y_t = Hy_{t,1:5}$.

1231 1232 **Data generation: initial draw of the parameters, continued.** The matrix A is drawn component-wise, in an i.i.d. manner: 1233 the $A_{i,j}$, where $i, j \in [5]$, follow a standard Gaussian distribution $\mathcal{N}(0, 1)$.

¹²³⁴ We draw $f = (f_1, \ldots, f_5)$ component-wise. To do so, we consider the following base functions $\mathbb{R}^3 \to \mathbb{R}$:

1241 We now explain how f_i is drawn for each component $i \in [5]$. First, the number k_i of effects to consider is drawn uniformly in 1242 the set [11]. Then, we sample with replacement k_i base functions in the set $\{g_1, \ldots, g_{11}\}$; we denote them by $h_{i,1}, \ldots, h_{i,k_i}$. 1243 Finally, we add signs: we draw k_i i.i.d. symmetric Rademacher random variables $r_{i,1}, \ldots, r_{i,k_i}$ (i.e., variables that take 1244 values -1 and 1 with respective probabilities 1/2). All in all, we let

$$f_i = \sum_{j=1}^{k_i} r_{i,j} h_{i,j}$$

1249 F.1.2. DATA SPLITTING

1212 1213

1214 1215

1216

1220 1221

1222

1226

1228

1250

1259 1260

We take $T = 100\,000$ (to contrast with the experiments on real data). These T observations are first randomly split in two subsets, containing 80% and 20% of the data.

The smaller subset is referred to as the test set and is denoted by $\mathcal{D}_{\text{test}}$. Its data points will play the role of the (x_{T+1}, y_{T+1}) , as explained later in Appendix F.1.4.

The larger subset of 80% of the data is split again in three sub-subsets, containing 40% (train set \mathcal{D}_{train}), 20% (estimation set \mathcal{D}_{estim}), and 20% (calibration set \mathcal{D}_{calib}) of the total data. These data points are used to construct the prediction rectangles, which are all (in Algorithms 1–2–3) of the form

$$oldsymbol{x} \longmapsto \prod_{i=1}^{8} \widetilde{C}_i(oldsymbol{x}) = \prod_{i=1}^{8} \left[\widetilde{\mu}_i(oldsymbol{x}) + \widetilde{q}_{lpha/2}^{(i)}, \, \widetilde{\mu}_i(oldsymbol{x}) + \widetilde{q}_{1-lpha/2}^{(i)}
ight]$$

and only depend on the features x through the centers $\tilde{\mu}_i(x)$. The algorithms that do not use an estimation set $\mathcal{D}_{\text{estim}}$, i.e., Algorithms 1–2, simply ignore data points in $\mathcal{D}_{\text{estim}}$. 1265 F.1.3. TRAIN SET: REGRESSION ALGORITHM \mathcal{A}

1284

1285 1286 1287

1292 1293 1294

1298 1299

1300 1301

The last piece to fully define the procedures implemented is to describe the regression algorithm \mathcal{A} given as input to Algorithms 1–2–3. This algorithm will be given by a base forecasting method run independently at each node.

Before we describe this base forecasting method, we mention a constraint that we impose. It turns out that in the practice of hierarchical forecasting, explanatory variables are not necessarily all available at every level of granularity within the hierarchical structure. (For instance, some meteorological data may only be available at specific locations equipped with the

1272 necessary sensors and cannot be communicated in a timely, real-time, manner to other nodes.) This also actually makes the 1273 hierarchy more interesting from a forecasting viewpoint since the observations at some nodes are harder to predict than 1274 others.

¹²⁷⁵ ¹²⁷⁶To reproduce this specificity, for each of the nodes at the most disaggregated level, indexed by $i \in [5]$, we draw independently ¹²⁷⁷a Bernoulli variable ρ_i with parameter 0.7: if $\rho_i = 1$, then the forecasting strategy may use the entire vectors \boldsymbol{x}_t ; otherwise, ¹²⁷⁸the forecasting strategy only accesses to $\boldsymbol{x}'_t = (x_{t,1}, x_{t,2})^{\mathsf{T}}$.

1279 It only remains to describe the forecasting strategy used independently at each node $i \in [8]$, based on features that lie in 1280 \mathbb{R}^2 or \mathbb{R}^3 . Given the additive nature (18) of the data, a natural choice is to resort to the theory of estimation of generalized 1281 additive models, see a reminder at the end of this subsection.

1282 1283 For each $i \in [8]$, depending on ρ_i , the regression estimate $\hat{\mu}_i$ produced for the *i*-th component of the *y* is of the form

 $\widehat{\mu}_{i}: \boldsymbol{x} \longmapsto \begin{cases} \widehat{\mu_{i}}^{(1)}(x_{1}) + \widehat{\mu_{i}}^{(2)}(x_{2}) + \widehat{\mu_{i}}^{(3)}(x_{3}), & \text{if } \rho_{i} = 1\\ \widehat{\mu_{i}}^{(1)}(x_{1}) + \widehat{\mu_{i}}^{(2)}(x_{2}), & \text{otherwise.} \end{cases}$

1288 **Reminder on generalized additive models.** Generalized additive models (GAMs, Wood, 2017) are a popular modeling 1289 for electricity demand. They form a good compromise between forecast efficiency and interpretability. In that setting, 1290 univariate response variables z_t based on features $x_t \in \mathbb{R}^d$, where $t \in [T]$, are expressed as

$$z_t = \beta_0 + \sum_{j=1}^d m_j(x_{t,j}) + \varepsilon_t , \qquad (19)$$

where the $m_j : \mathbb{R} \to \mathbb{R}$ do not depend on t and are called the non-linear effects, and where the ε_t are i.i.d. random noises. The non-linear effects m_j are each possibly decomposed on a given spline basis $(B_{j,k})$, chosen by the forecasting agent:

$$m_j: x \in \mathbb{R} \longmapsto \sum_{k=1}^{K_j} \beta_{j,k} B_{j,k}(x)$$

1302 where K_j depends on the dimension of the spline basis. Estimating the model (19) then amounts to estimating the 1303 coefficients $\beta_{j,k}$.

At a high level, we may write that the estimation of these coefficients $\beta_{j,k}$ is performed via by penalized least-squares, where the penalty term therein involves the second derivatives of the functions m_j , forcing the effects to be smooth. We resorted to the R package mgcv of Wood (2023) in our simulations, with the basis by default: the thin plate spline basis, with a maximum number of degrees of freedom of 10.

1309 F.1.4. Test set: evaluation of the prediction sets 1310

1311 The objectives (*) and (**) are both in terms of coverage probability and expected length with respect to all data (the 1312 observations to be predicted as well as the data used to compute the prediction intervals). We consider, for the expected-length 1313 criterion (**), weights given by the constant vector $w = (1, ..., 1)^{T}$.

On the test set, we estimate the conditional coverage probabilities and expected lengths given the specifications of the experiment (i.e., f, A, and the ρ_i) and given data in the sets $\mathcal{D}_{\text{train}}$, $\mathcal{D}_{\text{estim}}$, and $\mathcal{D}_{\text{calib}}$: for each $i \in [8]$,

$$c_i \stackrel{\text{def}}{=} \frac{1}{T_{\text{test}}} \sum_{t \in \mathcal{D}_{\text{test}}} \mathbb{1}_{\left\{y_{t,i} \in \widetilde{C}_i(\boldsymbol{x}_t)\right\}} \quad \text{and} \quad \ell_i \stackrel{\text{def}}{=} \ell\left(\widetilde{C}_i(\cdot)\right)$$

1320 where we denoted by $T_{\text{test}} = 20\,000$ the cardinality of $\mathcal{D}_{\text{test}}$. That is, for the conditional coverage probability, we resort to 1321 Monte-Carlo-type estimates. For the lengths of the intervals $\tilde{C}_i(\boldsymbol{x})$, we note that they do not depend on \boldsymbol{x} , so are constant on 1322 $\mathcal{D}_{\text{test}}$; we denote by $\ell(\tilde{C}_i(\cdot))$ their common value.

We actually run the entire procedure a large number of times to get unconditional probabilities and expectations, as described
 next.

1326 1327 F.1.5. Monte-Carlo estimates based on large numbers of runs

We run 1 000 the entire procedure and get, for each run, estimates of the conditional coverage probabilities and expected
 lengths, which we denote by:

1330 1331 1332

1335 1336

1340 1341 1342

1353 1354

1366 1367 1368 $c_i^{(r)} \quad \text{and} \quad \ell_i^{(r)} \,, \qquad \text{where} \quad i \in [8] \; \text{ and } \; r \in [1 \; 000] \,.$

We in turn get the following estimates for the unconditional coverage probabilities and expected squared lengths: for each $i \in [8]$, $i \in [8]$,

$$\bar{c}_i \stackrel{\text{def}}{=} \frac{1}{1\,000} \sum_{r=1}^{1\,000} c_i^{(r)} \quad \text{and} \quad \bar{\ell}_i \stackrel{\text{def}}{=} \frac{1}{1\,000} \sum_{r=1}^{1\,000} \left(\ell_i^{(r)}\right)^2.$$

These empirical means estimate the underlying unconditional coverage probabilities and expected squared lengths up to 95%–confidence errors margins given by

$$\gamma_{c,i} \stackrel{\text{def}}{=} 1.96 \frac{\text{std}\left(c_i^{(1)}, \dots, c_i^{(1\,000)}\right)}{\sqrt{1\,000}} \quad \text{and} \quad \gamma_{\ell,i} \stackrel{\text{def}}{=} 1.96 \frac{\text{std}\left(\left(\ell_i^{(1)}\right)^2, \dots, \left(\ell_i^{(1\,000)}\right)^2\right)}{\sqrt{1\,000}},$$

where std $(x_1, \ldots, x_{1\,000})$ denotes the standard deviation of the data series given as argument:

std
$$(x_1, \dots, x_{1\,000}) = \sqrt{\frac{1}{1\,000} \sum_{r=1}^{1\,000} (x_r - \overline{x}_{1\,000})^2}, \quad \text{where} \quad \overline{x}_{1\,000} = \frac{1}{1\,000} \sum_{r=1}^{1\,000} x_r.$$

For scaling issues on the lengths, we rather report, in our experiments, when dealing with component-wise quantities, the following point estimates and associated confidence intervals on the underlying unconditional probabilities and expectations: for all $i \in [8]$,

$$\overline{c}_i \quad \text{and} \quad \sqrt{\overline{\ell}_i}, \qquad [\overline{c}_i \pm \gamma_{c,i}] \quad \text{and} \quad \left[\sqrt{\overline{\ell}_i - \gamma_{\ell,i}}, \sqrt{\overline{\ell}_i + \gamma_{\ell,i}}\right].$$
 (20)

1355 F.1.6. COMPONENT-WISE RESULTS: COVERAGE AND LENGTH

The top graph of Figure 4 reports the indicators defined in (20), with standard errors in x-axis corresponding to the estimation of the component-wise coverage levels, and standard errors in y-axis, to the ones for the lengths. The comments on the outcomes are to be found in Section 4.

1360 1361 F.1.7. Global results: total lengths

We also report results on the total lengths, i.e., for the quantities appearing in the efficiency objective (**), where we recall that w = 1.

In the same spirit as in Section F.1.5, we consider 1365

$$L_{\bullet}^{(r)} = \sum_{i=1}^{8} \left(\ell_{i}^{(r)}\right)^{2}, \quad \text{where } r \in [1\,000], \qquad \overline{L}_{\bullet} \stackrel{\text{def}}{=} \frac{1}{1\,000} \sum_{r=1}^{1\,000} L_{\bullet}^{(r)} = \sum_{i=1}^{8} \overline{\ell}_{i}.$$

This empirical mean estimates the underlying expected sum of the squared lengths up to 95%-confidence errors margins given by

- $\gamma_{L,\bullet} \stackrel{\text{def}}{=} 1.96 \frac{\text{std}\left(L_{\bullet}^{(1)}, \dots, L_{\bullet}^{(1\ 000)}\right)}{\sqrt{1\ 000}},$ 1374
 - 25



Figure 4. Artificially generated data: component-wise coverage levels and prediction-interval lengths (*top figure*) and total lengths (*bottom figure*). This figure merely performs a zoom on the left graphs of Figure 2. The standard errors reported are based on the formulas (20) and (21).

1430 F.2. Palo Alto daily charging Energy

1431 The dataset we consider is presented in greater detail by Amara-Ouali et al. (2021), an article referencing several data sets 1432 for charging sessions of electric vehicles. We refer to this data set as the Palo Alto dataset, as it is related to charging stations 1433 located in the city of Palo Alto, CA. It contains a substantial number of charging sessions and an interesting hierarchical 1434 structure, with 47 charging points that are divided into a dozen of stations. However, due to some real-world considerations, 1435 data is not available for all these charging points in the 2015–2019 period considered. We only consider 2 charging stations 1436 (called "Riconada Library" and "Hamilton", featuring 3 and 2 charging points, respectively) for which data is available 1437 on the entire period. The study stops in 2019 to avoid the temporary shift in demand caused by Covid19 in 2020 and the 1438 subsequent years. 1439

We are interested in daily energy demand at each charging point, each charging station, and at the global level of the considered hierarchy. The total number of observations for each node is T = 1780.

1443 F.2.1. METHODOLOGY

1477

1482

1484

The methodology followed for this data set essentially mimics the one for artificial data, as presented in Section F.1. We thus present only the main adaptations made.

Runs and data splitting. We will perform 360 runs of a given experimental procedure, described next. We split into train, estimation, calibration, and test sets as in Section F.1.2, with same 40%-20%-20%-20% proportions.

Regression algorithm \mathcal{A} . As in Section F.1.3 (and taking inspiration from Amara-Ouali et al., 2022), we resort to modeling and forecasting through GAMs. More specifically, for each node $i \in [8]$ and day t, we consider the following auto-regressive specification of GAM for the energy demand $y_{t,i}$:

$$y_{t,i} = \beta^{(0)} + \sum_{j=1}^{7} \beta_j^{(1)} \mathbb{1}_{\text{DayType}_t = j} + \sum_{j=1}^{7} m_1(y_{t-1}) \mathbb{1}_{\text{DayType}_t = j}$$
(22)

$$+ m_2(y_{t-1,i}) + m_3(y_{t-7,i}) + m_4(t) + m_5(\text{To}\mathbf{Y}_t) + \varepsilon_t, \qquad (23)$$

where DayType_t $\in \{1, ..., 7\}$ is a categorical variable indicating the day of the week, ToY_t is the "time of year", i.e., the position of the day in the year (whose value grows linearly from 0 on the 1st of January to 1 on the 31st of December). The model (22) incorporates a trend term $m_4(t)$, which may be estimated because we pick a random subset of the entire data set for the train set; this trend term is useful to take into account changes in the infrastructures and shifts in user behaviors.

We again resorted the R package mgcv of Wood (2023) to forecast this model and get $\hat{\mu}$. The parameters were the default thin plate spline basis and a maximum number of degrees of freedom of 10 for the estimation of the coefficients for $m_1, m_2 m_3$, and 15 for m_4 . We fitted m_5 with cyclic splines and a maximum number of degrees of freedom of 30.

Validity check of Assumption 3.6 on elliptic distribution. Based on the regression function $\hat{\mu}$ output, we computed the signed non-conformity scores on the estimation and calibration sets. We then fitted a Student distribution (which is a particular case of an elliptic distribution). Figures 5 and 6 illustrate the goodness of fit between the empirical distributions of scores and the Student distributions with parameters estimated on these scores, through, respectively, densities and Q-Q-plots. We were interested in two nodes *i*: the total node (the total demand for the 5 charging points) and the node of the Riconada Library (the sum of the 3 charging points located there).

1474 The fit to a Student distribution looks reasonable in both cases, with actually an excellent fit for the total node and a relatively 1475 small mode of the distribution located in the right tail of the distribution scores being not captured well in the case of the 1476 Riconada Library node.

1478 Estimation set: details on the "Oracle MinT" strategy. For this data set, we also report the performance of the "Oracle 1479 MinT" strategy, i.e., Algorithm 2 used with $P_{\Sigma^{-1}}$. The question is how to determine Σ .

¹⁴⁸⁰ 1481 The "MinT" strategy, i.e., Algorithm 3 with \mathcal{P}_{MinT} , truly estimates Σ , on the estimation set:

$$\widehat{\Sigma} = \frac{1}{T_{\text{estim}}} \sum_{t \in \mathcal{D}_{\text{estim}}} \left(\widehat{s}_t - \overline{s} \right) \left(\widehat{s}_t - \overline{s} \right)^{\top}, \quad \text{where} \quad \overline{s} = \frac{1}{T_{\text{estim}}} \sum_{t \in \mathcal{D}_{\text{estim}}} \widehat{s}_t.$$





Figure 5. Student density estimator compared to the empirical density of the scores for the total node (a) and the Rinconada Library station (b).



Figure 6. Q-Q-plots of the scores for the total node (a) and the Rinconada Library station (b).



Figure 7. The estimates $\hat{\Sigma}$ and $\hat{\Sigma}^*$ of the covariance matrix of the non-conformity scores for one given run of the experiment.

1539

1519

1540 In the case of the Palo Alto data set, the results obtained by this strategy are rather poor, which may be attributed to having 1541 few data points only to perform the estimation: the ones of the estimation set \mathcal{D}_{estim} .

To determine Σ for the "Oracle MinT" strategy, we actually also estimate it, but by cheating: we produce an estimate thereof using \mathcal{D}_{calib} and \mathcal{D}_{test} , the two subsets on which the non-conformity scores will be calculated. More precisely, we produce the estimate

$$\widehat{\Sigma}^{\star} = \frac{1}{T_{\text{calib}} + T_{\text{test}}} \sum_{t \in \mathcal{D}_{\text{calib}} \cup \mathcal{D}_{\text{test}}} \left(\widehat{s}_t - \overline{s}^{\star} \right) \left(\widehat{s}_t - \overline{s}^{\star} \right)^{\top}, \quad \text{where} \quad \overline{s}^{\star} = \frac{1}{T_{\text{calib}} + T_{\text{test}}} \sum_{t \in \mathcal{D}_{\text{calib}} \cup \mathcal{D}_{\text{test}}} \widehat{s}_t ,$$

1549 and run Algorithm 2 with $P_{(\widehat{\Sigma}^{\star})^{-1}}$.

Figure 7 displays the two estimates $\hat{\Sigma}$ and $\hat{\Sigma}^{\star}$ just for one given run picked at random out of the 360 runs. The differences between the two estimates look mild, yet, the differences in performance are substantial (see Figure 8). We repeated the comparison several times and always obtained quite similar estimates.

1555 F.2.2. OUTCOMES

We use the exact same metrics as in Sections F.1.5 and F.1.7, with the mere replacement of the number 1 000 of runs considered therein by the number 360 of runs considered now, and report the corresponding results in Figure 8. This figure merely performs a zoom on the right graphs of Figure 2. The comments on the outcomes are located in Section 4.



Figure 8. Palo Alta data of daily energy demand for charging electric vehicles: component-wise coverage levels and prediction-interval lengths (*top figure*) and total lengths (*bottom figure*). This figure merely performs a zoom on the right graphs of Figure 2. The standard errors reported are based on the formulas (20) and (21), with 1 000 replaced by 360.